

1 **The Choice of Response Alternatives**
2 **in COVID-19 Social Science Surveys**

3
4 Daniel B. Wright¹, Sarah M. Wolff¹, Rusi Jaspal²,
5 Julie Barnett³, and Glynis M. Breakwell^{3,4}

6 ¹University of Nevada, Las Vegas

7 ²University of Brighton

8 ³University of Bath

9 ⁴Imperial College, London

10 **Author Note**

11 The research is support by British Academy award CRUSA210025 to JR, GB, JB,
12 and DW, and this provided funding for the first study. DW and SW are supported by an
13 endowment from the Dunn Family Foundation and this provided funding for the second
14 study. The second study was pre-registered at OSF: <https://osf.io/b4uef/> and received
15 IRB approval from the UNLV IRB [1753484-2]. The authors have no conflicts of interest.
16 The Imperial Data set is available at
17 <https://github.com/YouGov-Data/covid-19-tracker> (Jones, 2020). The de-identified
18 data file for the second study are at:
19 <https://github.com/dbrookswr/RespAlt/blob/main/WWRespAlt.csv>. The code for all
20 the analyses will appear as the final version as a reproducible **knitr** (Y. Xie, 2013) document
21 on [github](https://github.com). Correspondence should be sent to [daniel.wright](mailto:daniel.wright@unlv.edu) at unlv.edu.

Abstract

22

23 Social science research is key for understanding and for predicting compliance with
24 COVID-19 guidelines, and much of this research relies on survey data. While much focus is
25 on the survey question stems, less is on the response alternatives presented that both
26 constrain responses and convey information about the assumed expectations of the survey
27 designers. The focus here is on the choice of response alternatives for the types of behavioral
28 frequency questions used in many COVID-19 and other health surveys. We examine issues
29 with two types of response alternatives. The first are vague quantifiers, like “rarely” and
30 “frequently.” Using data from 30 countries from the Imperial COVID data hub, we show
31 that the interpretation of these vague quantifiers (and their translations) depends on the
32 norms in that country. If the mean amount of hand washing in your country is high, it is
33 likely “frequently” corresponds to a higher numeric value for hand washing than if the mean
34 in your country is low. The second type are precise numeric response alternatives and they
35 can also be problematic. Using a US survey, respondents were randomly allocated to receive
36 either response alternatives where most of the scale corresponds to low frequencies or where
37 most of the scale corresponds to high frequencies. Those given the low frequency set
38 provided lower estimates of the health behaviors. The choice of response alternatives for
39 behavioral frequency questions can affect the estimates of health behaviors. How the
40 response alternatives mold the responses should be taken into account for epidemiological
41 modeling. We conclude with some recommendations for response alternatives for health
42 behavioral frequency questions in surveys.

The Choice of Response Alternatives in COVID-19 Social Science Surveys

People's reactions to mask and COVID-19 vaccine regulations lay bare the need for psychological research to complement biological and economic research for effective management of epidemics. The underlying data for much psychological research come from responses to sample surveys. Constructing survey questions that lead to valid, reliable, and fair responses is difficult (Groves et al., 2009). In everyday conversations, people ask questions without constraining others to respond using pre-defined formats with limited options, but in survey conversations this often occurs. This is done to ease the coding of responses and often to get the respondents to translate their complex beliefs into a single value that is more suitable for statistical analyses.

The responses from health surveys are critical for monitoring health related behaviors, evaluating health campaigns, understanding/modeling the spread of disease, and developing public policy (e.g., Gadarian, Goodman, & Pepinsky, 2021). Health behavior data inform resource allocation decisions and provide necessary information to identify target groups for intervention, to track progress, and to evaluate existing strategies (e.g., Lee & Thacker, 2011). The current research was prompted by the COVID-19 epidemic and the realization that one of the main reasons for its level of impact in many countries is people failing to heed guidelines from scientific groups on the efficacy of health related behaviors (e.g., hand washing, mask wearing, vaccines). Epidemiologists use estimates from social surveys to gauge how much people follow these guidelines.

Self-reports are prone to bias and memory errors (Sudman & Bradburn, 1973). In order to understand measurement error within surveys, researchers examine the cognitive processes that occur when respondents answer questions (e.g., Belli, Conrad, & Wright, 2007; Fienberg, Loftus, & Tanur, 1985; Loftus, Fienberg, & Tanur, 1985; Schaeffer & Dykema, 2020; Tourangeau, 2003). The theoretical approach taken here is to assume the survey

70 situation is an artificial conversation, and like other conversations respondents use the
71 information presented to them to interpret the meaning of questions (Schwarz, 1995). For
72 modeling epidemiology, the focus is often on estimates from behavioral frequency questions
73 (Burton & Blair, 1991), like “how often do you wash your hands?” Our focus is on the
74 response alternatives. As these are not part of most everyday conversations, when they are
75 presented in surveys they may stand out. When responding to a survey question, the
76 responses are often ordered, and respondents will see this as a scale composed of words.
77 According to Grice’s maxims (1975, see also Wilson and Sperber, 2012) when people are
78 presented with a question they assume that the information--including that within the
79 prescribed response format--will be accurate and relevant. While the survey is an artificial
80 conversation, respondents may still assume that the scale and words are appropriate for the
81 behavior in the question (e.g., Schwarz, 1995).

82 Respondents assume that researchers construct a meaningful scale that reflects
83 appropriate knowledge about the distribution of the behavior. Accordingly,
84 values in the middle range of the scale are assumed to reflect the ‘average’ or
85 ‘typical’ behavior, whereas the extremes of the scale are assumed to correspond
86 to the extremes of the distribution. Schwarz (2010, p. 49)

87 If respondents believe the information implied by the set of response alternatives is at
88 odds with what they believe, they could assume the response alternatives were not chosen to
89 have the property Schwarz describes. This could occur for many surveys because some online
90 surveys are poor quality and sometimes designers use the same scale for multiple behaviors
91 without concern of whether it is appropriate for each of these behaviors. This lessens the
92 information likely gathered from the survey and may lower respondents’ view of the
93 designers. Alternatively, the scale could affect what respondents think about the behavior
94 and how they respond (see Figure 1). For example, it may change what respondents think
95 about the behavior in the question stem. If a person is unsure of population norms and is
96 presented with a set of response alternatives suggesting the behavior is common, they may

97 come to believe that the behavior is common. Another possibility is that the response
98 alternatives could change what the respondents think the target event is. Wright, Gaskell,
99 and O’Muircheartaigh (1997) describe how this can occur for vague and ambiguous terms.
100 For example, they asked respondents how often their teeth were cleaned either with response
101 alternatives suggesting this meant by a dentist or with response alternatives suggesting this
102 meant by themselves. Which set respondents received affected what respondents thought
103 the target behavior was. This also occurs for vague behaviors like being annoyed or being
104 satisfied (Gaskell, O’Muircheartaigh, & Wright, 1994; Schwarz, Strack, Müller, & Chassein,
105 1988). For relatively well defined events (e.g., how many cups of coffee you have in a typical
106 day), the definitions should not be greatly affected. If the difference between the population
107 norms implied by the response alternatives and their pre-survey beliefs is great, this might
108 cause them to doubt the applicability of the scale, or to change their normative beliefs.
109 Finally, respondents can interpret any set of alternatives as a scale from low to high,
110 ignoring the particular words used to compose the scale. This will be more likely when the
111 response alternatives are vague. In these cases it is unclear what question they answer: how
112 much they engage in the behavior compared with others; with their expectations; with their
113 behavior before the pandemic; etc. This is discussed further at the end of the paper when we
114 make recommendations for the choice of response alternatives.

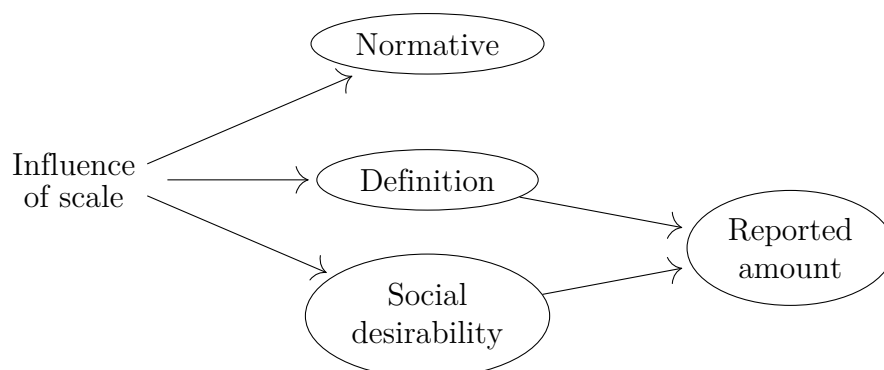


Figure 1

How the response alternatives can affect respondents’ beliefs about the meaning of the terms, about population norms, and their responses.

115 The response alternatives for behavioral frequency questions generally have one of
116 five formats (more elaborate approaches exist, for example having respondents list the
117 behaviors on a calendar and some of these are discussed in the recommendations section, e.g.,
118 Schatz, Knight, Belli, & Mojola, 2020), each with its own concerns:

- 119 1. Free recall. Respondents provide a numerical estimate, often for a specific duration.
120 This can make difficult memory demands for high frequency behaviors. Respondents
121 often use rough heuristics to *gestimate* the frequency of the behavior. If they try to
122 recall every incident this tends to lead to under-reporting. Another issue is that
123 respondents often round their responses, giving response prototypes (Huttenlocher,
124 Hedges, & Duncan, 1991).
- 125 2. Last time. Respondents can be asked the last time a relatively infrequent event
126 happened, or if it happened since some memorable event (e.g., Loftus & Marburger,
127 1983). Time since an event can be used to estimate the event frequency. The difficulty
128 is people have difficulty remembering *when* events occur and tend to forward telescope
129 these dates to be more recent than is accurate (e.g., Neter & Waksberg, 1964;
130 Thompson, Skowronski, & Lee, 1988).
- 131 3. Comparison. Respondents say how often they experience the behavior compared with
132 others. The difficulty is respondents may not know how often others experience the
133 behavior.
- 134 4. Vague quantifiers. The respondents use a set of alternatives composed of vague
135 quantifiers to describe how much they experience the behavior. Study 1 examines how
136 different people interpret these differently.
- 137 5. Numeric response alternatives. Respondents can be provided with a set of exhaustive
138 and mutually exclusive numeric response categories. The concern here, examined in
139 Study 2, is how the choice of these sets can affect how the respondent answers the
140 question and therefore the study's results.

141 Before discussing the empirical research, it is important to stress that one of the main
 142 limitations of behavioral frequency question is the inaccuracies of human memory. While
 143 this is not addressed here, it is important that survey designers take this and other cognitive
 144 fallibilities into account when creating surveys.

145 **Study 1: Vague Quantifiers**

146 Schaeffer (1991) titled her paper “Hardly ever or constantly” as reference to a scene
 147 in the film *Annie Hall* where Alvie Singer and Annie Hall are each ask how often they sleep
 148 together:

149 Alvie: Hardly ever, maybe three times a week.
Annie: Constantly, I’d say three times a week.

150 This highlights that vague quantifiers can be associated with different numeric values for
 151 different people.

152 Wright, Gaskell, and O’Muircheartaigh (1994) showed that part of the differences in
 153 how people used vague quantifiers can be attributed to what they think the normative
 154 behaviors of that behavior are. They asked UK respondents how much they thought people
 155 typical watched television, and found large differences in this belief by social class. In a
 156 subsequent study they asked the following two questions.

157	Q1 And from this list, on average, how much television do you watch on a typical weekday?	None at all Hardly any A little Quite a bit A lot
	Q2 And about how many hours would that be?	[Free Recall]

158 They found that respondents from classes that watched more television and thought the
 159 normative behavior was higher believed each vague quantifier corresponded to more hours
 160 than did respondents from classes that watched less television. Groups that thought people
 161 tended to watch more television interpreted the vague quantifiers as corresponding to higher
 162 amounts. Using a large multi-country database, we examine the association between

163 responses using vague quantifiers and a numeric response. Both of these attempt to measure
164 the frequency of the behavior.

165 Rather than social class, the current study estimates the amount each country does
166 the target behavior and examines if this is a good indicator of how people in the country
167 interpret the vague quantifiers. This study uses data from thirty countries and the survey
168 was delivered in many languages. Linguistic differences make cross-country, and within
169 country where multiple languages are used, difficult. Our prediction is that countries where
170 there is a higher mean for the numeric response, will have higher numeric values
171 corresponding to the different vague quantifiers (and their translations). We split the data
172 into 25% for estimating the means for the behaviors and then used the remaining 75% to
173 test the prediction.

174 **Methods**

175 Data from <https://github.com/YouGov-Data/covid-19-tracker> (Jones, 2020)
176 were downloaded on November 8, 2021. The database and discussion of their methods are
177 available at www.coviddatahub.com. In total, at the time of download, there were 734,075
178 respondents. There are five questions of interest. The first variable is the country. There are
179 thirty countries. Sample sizes of those with complete responses are shown in Table 1.

180 The next three variables of interest are two vague quantifier questions and one free
181 recall question about hand-washing/sanitizing. These were separated by 17 questions. Their
182 wording, from the UK version (i.e., *sanitising*), is shown in Table 2. Only data where there
183 were no missing values for these were used, leaving 646,177 cases. The final question is the
184 Cantril ladder (1965), which asks respondents to imagine a ladder from 0 to 10 steps and
185 asks them to rate their current life satisfaction. About 9% have missing values for this. It is
186 used for exploratory purposes at the end of this section.

187 **Results and Discussion**

188 The vague quantifier questions are not directly comparable to the free recall question
189 as the latter combines the two behaviors asked about in the vague quantifier questions. We

Table 1

Countries in the Imperial database, number of respondents with complete data for the questions analyzed here, and the mean for the natural logarithm of the responses to the free recall question, plus a starting value of +.5, for the 25% of the training sample (given the sample sizes these are very similar to the total means).

Country	n	mean(ln(x + .5))	Country	n	mean(ln(x + .5))
Australia	35,932	2.03	Brazil	10,308	2.26
Canada	31,060	2.15	China	15,985	1.30
Denmark	32,154	2.39	Emirates	11,924	2.16
Finland	19,076	2.07	France	36,360	2.07
Germany	36,271	1.93	Hong Kong	7,042	1.77
India	16,145	2.07	Indonesia	12,141	1.90
Israel	6,480	1.88	Italy	36,106	2.19
Japan	16,931	1.49	Malaysia	12,133	1.91
Mexico	12,012	2.31	Netherlands	11,399	1.93
Norway	32,032	2.21	Philippines	12,002	2.07
Saudi Arabia	11,450	1.95	Singapore	32,977	1.84
South Korea	15,644	1.72	Spain	36,198	2.06
Sweden	36,201	2.11	Taiwan	11,993	1.53
Thailand	12,085	1.76	UK	46,143	2.16
US	27,919	2.02	Vietnam	12,074	1.75

190 begin by comparing responses from each of vague quantifier questions with the free recall
 191 responses question. A small percentage of respondents (0.07%) gave responses of 1,000 or
 192 more to the free recall question. Assuming these respondents are awake for 18 hours, this is
 193 about once a minute. This is a very skewed variable (skew = 25.06). It was transformed
 194 using $\ln(x + .5)$ (the +.5 as some people, 1.22%, said zero) and this lessened the skewness to
 195 0.14.

196 The associations between the vague quantifier responses and the transformed numeric
 197 responses from the free recall question are shown in Figure 2. The relationship between these
 198 two variables are fairly weak, even allowing for them asking asking about slightly different
 199 behaviors. Twenty-five percent of the data were used to estimate the mean for each country
 200 for the transformed ($\ln(x + .5)$) responses from the free recall question. These are shown in
 201 Table 1. The remaining 75% of the data are used to explore the relationship between the

Table 2

Vague quantifier and free recall hand washing/santizing questions from the Imperial data set.

Question	Resp. Alts.
Washed hands with soap and water	Always Frequently Sometimes Rarely Not at all
Used hand sanitiser	Always Frequently Sometimes Rarely Not at all
Thinking about yesterday ... about how many times, would you say you washed your hands with soap or used hand sanitiser?	[Free recall]

Variable names are as found in the Imperial webpage.

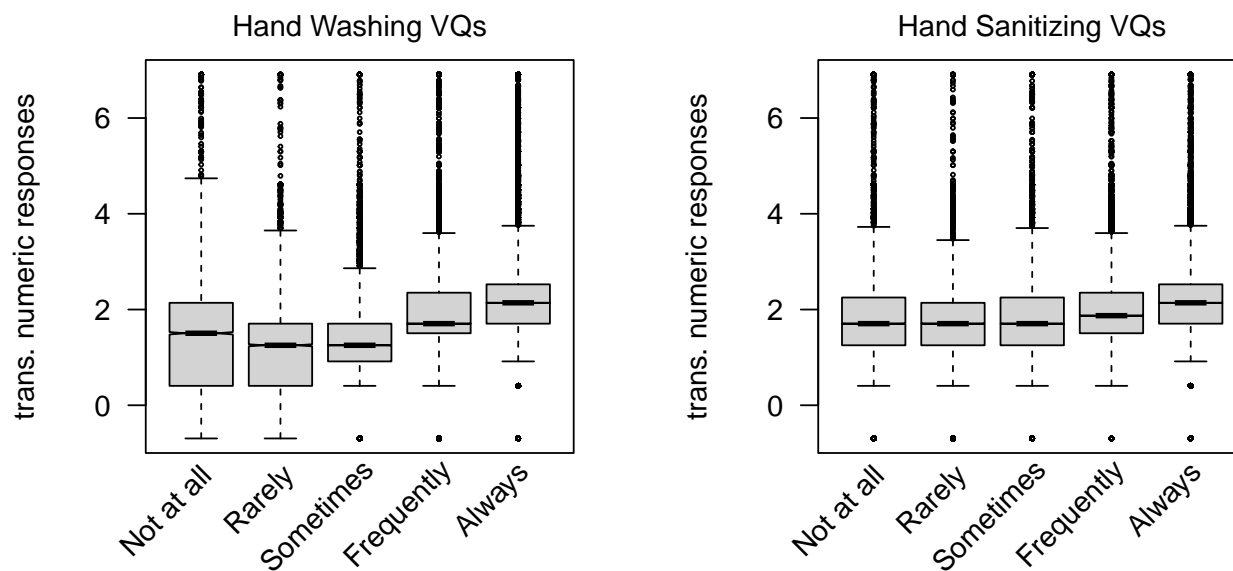


Figure 2

Boxplots for the transformed numeric response from the free recall question for the two vague quantifier questions.

202 transformed numeric responses and the vague quantifiers. Given the large sample size, the
203 25% and 75% are sufficient for our purposes. Since there are separate vague quantifier
204 questions for hand washing and hand sanitizing, these will be examined separately.

205 The hand washing vague quantifier variable is treated as categorical, so with $df = 4$
206 for its five categories. When it is used to predict the transformed variable from the free
207 recall question, the R^2 value was .119. Including the categorical variable of country, with its
208 $df = 29$, increased this to .177. The critical question is how much of this difference,
209 $\Delta R^2 = .058$, can be accounted for by the single variable ($df=1$) corresponding to the mean of
210 the transformed variable taken from the other 25% of the total sample? If there was 1/29th
211 increase the value would be about 3% of this amount, so approximately: .121. It was
212 $R^2 = .175$, or 96.38% of the possible amount. Figure 3 shows this. The color corresponds to
213 the mean for the transformed free recall responses. The greener the line the lower the
214 country mean for the transformed values (the color is based on a gradient between green and
215 red, so countries with means near the middle appear brown-ish). As is clear, the greener
216 lines are lower than the redder lines.

217 The findings were similar comparing the hand sanitizing vague quantifier variable
218 with the transformed free recall responses. When just the vague quantifier variable is used to
219 predict the these the R^2 value was .081. Including the categorical variable of country
220 increased this to .143, and difference of $\Delta R^2 = .058$. Using the single country mean variable
221 produced an $R^2 = .140$, or 93.97% of the possible amount (as opposed to 3%). This is shown
222 in Figure 4 with the greener lines being below the redder lines.

223 In Study 2 we examine the relationships between response alternatives effects and
224 some attitude questions. Therefore we felt it prudent to examine if an attitude variable,
225 self-expressed life satisfaction (the Cantril Ladder) mediated the relationship between the
226 vague quantifiers and the free recall variable. The Cantril ladder is weakly associated with
227 the transformed free recall responses: Pearson's $r = .072$. Our interest was whether it has
228 predictive value after accounted for the vague quantifier variables. The R^2 for predicting the

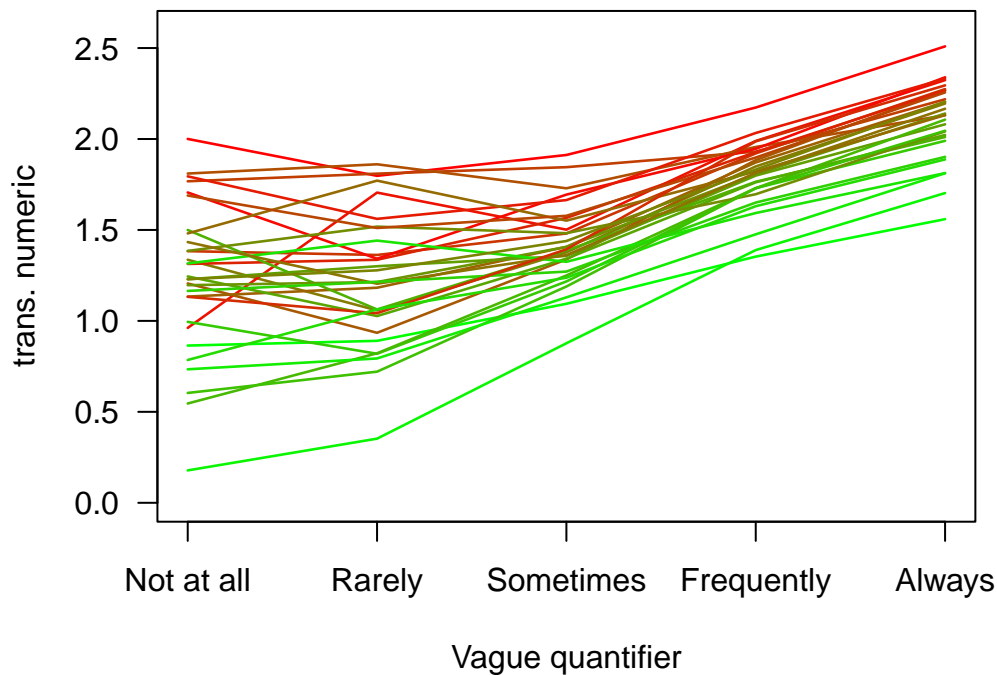


Figure 3

*The relationship between the means of the transformed free recall responses for each vague quantifier by country for the vague quantifier washing question. Countries with low free recall responses from the 25% of the sample are shown in **green** and those with high values in **red**, with intermediary countries shown in a mixture of these two colors.*

229 transformed free recall variable using the hand washing vague quantifier was $R^2 = .123$ and
 230 using the hand sanitizing vague quantifier was $R^2 = .089$. Including the interactions between
 231 Cantril's ladder and each vague quantifier question raised these to $R^2 = .124$ and to
 232 $R^2 = .090$, respectively. These increases are small and will not be considered further.

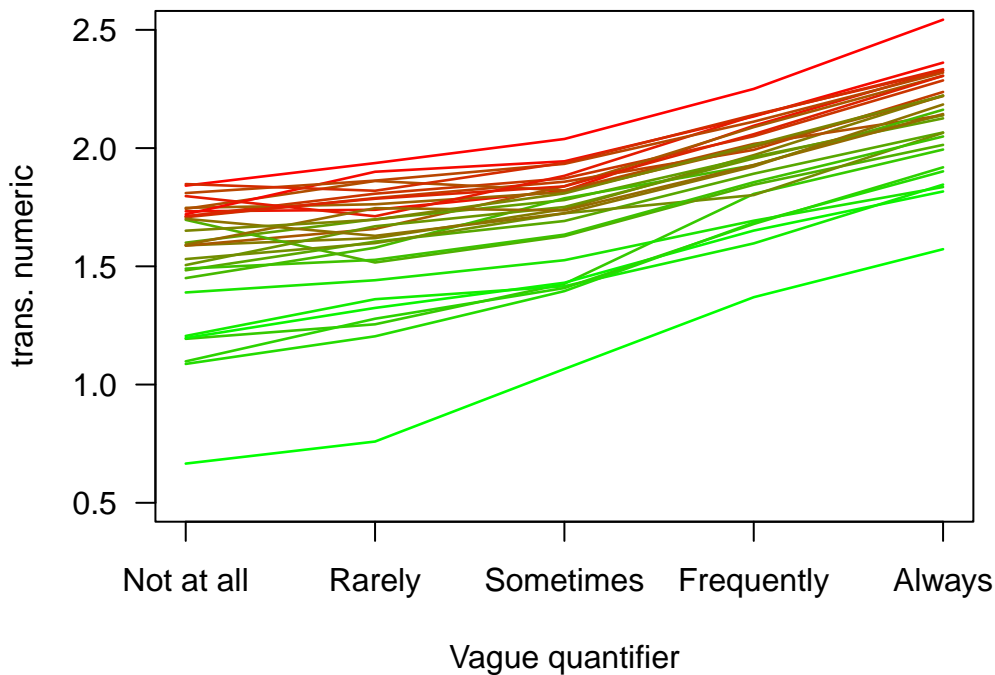


Figure 4
The relationship between the means of the transformed numeric estimates for each vague quantifier by country for the vague quantifier hand sanitizer question. Countries with low free recall responses from the 25% of the sample are shown in green and those with high values in red, with intermediary countries shown in a mixture of these two colors.

233

Study 2: Numeric Response Alternatives

234

The purpose of the COVID-19 Household Pulse Survey was to examine the effects of the pandemic on variables ranging from mental health behaviors to vaccine perceptions (US Census Bureau, 2021). The data are meant to help government direct aid, assistance, and support to the people and places that need it most. An example question reads:

237

<p>238</p> <p>Over the last 7 days, how often have you been bothered by the following problems ... Feeling nervous, anxious, or on edge? Would you say not at all, several days, more than half the days, or nearly every day? <i>Select only one answer</i></p>	<p>Not at all</p> <p>Several days</p> <p>More than half the days</p> <p>Nearly every day</p>
---	--

239 Schwarz and colleagues (Schwarz, 1995; Schwarz & Hippler, 1987; Schwarz, Hippler,
240 Deutsch, & Strack, 1985; Schwarz et al., 1988) conducted a series of studies to show how the
241 choice of which response alternatives to present can affect behavior estimates (see also
242 Gaskell et al., 1994, for use in national surveys). Schwarz uses Grice's maxims of
243 communication (Grice, 1975) applied to the survey situation explain his findings. Consider
244 the survey question above. The question asks respondents to reflect over the last several
245 days. Three of the four response alternatives involve the event happening multiple days in
246 the previous week: several days; more than half the days; and nearly every day. According
247 to Schwarz this may make respondents feel that being nervous, anxious, or on edge, are
248 likely to occur with greater frequency in the population than if the response alternatives
249 were: more than once a month, once a month, and none.

250 The behaviors used by Schwarz, Gaskell, and their colleagues were specifically chosen
251 to show the survey methodological effects predicted by Schwarz's hypotheses. The behaviors
252 used here were chosen because of their relation to disease transmission and that they are
253 part of health guidelines (e.g., those from the Center for Disease Control and Prevention
254 [CDC] in the US, and the National Health Service [NHS] in the UK). Also, Schwarz et al.
255 used in person (face-to-face), pen-and-paper, and telephone administration modes.
256 Nowadays, online surveys are becoming more common so the current study uses the online
257 administration mode. Our primary research question is whether the choice of response
258 alternatives affects the estimates of several health related behaviors. Respondents are
259 randomly allocated into one of two conditions. Those in the first condition are asked
260 behavioral frequency questions with response alternatives that discriminate more finely at
261 low frequencies and those in the second condition are asked these questions with response
262 alternatives that discriminate more finely at high frequencies.

263 **Methods**

264 The study received IRB approval from the UNLV IRB [1753484-2]. The authors have
265 no conflicts of interest. This study was pre-registered at <https://osf.io/b4uef/> and the

266 data are available at

267 <https://github.com/dbrookswr/RespAlt/blob/main/WWRespAlt.csv>

268 ***Sample***

269 Respondents were recruited via Amazon Mechanical Turk (MTurk). To be an MTurk
 270 worker you need to be 18 or over and with US social security number. Additional inclusion
 271 restrictions were that respondents had to have 100 previous what are called MTurk human
 272 intelligence tasks (HITs) and to have at least a 95% satisfaction rating from those conducting
 273 the research. Respondents were compensated \$2 for completing the questionnaire.

Table 3

Respondent flow showing the allocation into conditions and those excluded for a duplicate IP addresses or responding too quickly.

	Number of Respondents: Total n = 695	
	Low Response Alternatives	High Response Alternatives
Total Recruited	345	350
Duplicate IP	16	19
Too Quick	9	10
Total excluded	25	29
Percentage excluded	7%	8%
Final analysed (total n = 641)	320	321

274 A *catcha* question was included in the Qualtrics survey. Not correctly completing
 275 this meant the survey would not be included, but all passed this (or did not complete the
 276 survey). There were two additional exclusion criteria: multiple uses of an IP address and
 277 responding too quickly. While each MTurk worker needs a separate US social security
 278 number, people could use multiple MTurk accounts. A more likely reason is that people
 279 from within the same household are responding using the same IP address. As these people
 280 may talk about the study before the second completes it (and would in other ways also be
 281 non-independent), the duplicates (i.e., not the first one using the IP address) were excluded.
 282 Respondents who on average responded faster than two seconds for the behavioral frequency

283 and attitude questions were also excluded. The number of people excluded for these reasons,
284 in both conditions, is shown in Table 3.

285 *Materials and Experimental Design*

286 Respondents were asked three behavior/health questions related to COVID-19:

- 287 • How often did you wash your hands?
- 288 • When you washed your hands, typically how long did you spend?
- 289 • In a typical day, how often did you apply hand sanitizer?

290 They were randomly assigned to have either the low or the high response alternatives as
291 listed in Table 4.

292 Respondents were then asked seven attitude questions are about their health beliefs.
293 Respondents used a 0--100 sliding scale. Responses were measured to the tenth, for example
294 29.4. Our analyses concerning these variables are exploratory and concern whether the set of
295 response alternatives that were presented affects the responses on these variables.

- 296 • Compared with other people, how much were you concerned with the economic
297 impacts of the pandemic?
- 298 • Compared with other people, how much were you concerned with the health impacts
299 of the pandemic?
- 300 • Thinking back to the previous twelve months, how concerned were you about catching
301 COVID-19 yourself?
- 302 • Thinking back to the previous twelve months, how concerned were you about yours
303 friends and family catching COVID-19?
- 304 • Suppose that you were supposed to meet a small group of people. If you were not
305 feeling well (slight fever, running nose), how likely would you have stayed at home?
- 306 • Suppose that you were not feeling well (slight fever, running nose), how likely is it that
307 you would have consulted a medical professional?

Table 4

The response alternatives for the low and the high frequency conditions. The dashed lines show how the raw data can be re-coded for comparable frequencies/durations.

Question	Low	High
Hand washing frequency	Never	
	1	3 or less
	2--3	
	4--6	4--6
	7 or more	7--10
Hand wash duration		11--20
	less than 2--3 seconds	More than twenty
	5 seconds	5 seconds or less
	10 seconds	10 seconds
	20 seconds	20 seconds
Hand sanitiser	More than 20 seconds	One minute
		More than one minute
Hand sanitiser	Never	
	1	3 or less
	2--3	
	4--6	4--6
	7 or more	7--10
	11--20	More than twenty

- 308 • People vary in how much they trust scientists with respect to many issues. Please rate
309 your view.

310 In addition, respondents were asked their year of birth, gender, ethnicity, and asked to rate
311 their political beliefs on a 0--100 liberalism/conservatism scale.

312 The data were collected on a Monday evening, May 17, 2021, and this was a few days
313 after President Biden relayed CDC advice that masks need not be worn by fully vaccinated
314 people indoors in the USA (later, with the increased spread of new variants, guidance
315 changed).

316 *Statistical Plan*

317 The behavioral frequency questions are treated in two ways. When treating them as
318 1--5 rating scales, the means of these three are compared between the two groups using
319 Hotelling's T^2 . When the responses are re-coded into matching semantic categories (within
320 the dashed lines of Table 4), they are treated as categorical variables and compared using χ^2
321 tests, Cramér's V s, and multinomial logistic regressions. The associations among the
322 attitude questions are examined. The individual item distributions are skewed, so the data
323 are transformed. The associations suggest one underlying dimension. Scores on this
324 dimension are compared for the two groups.

325 **Results and Discussion**

326 *Behavioral Frequency/Duration Questions*

327 Three related sets of statistical analyses are conducted for examining the behavioral
328 frequency estimates in this study. The first examines if there are differences in responses for
329 the two groups if the response labels are not considered. As such, the questions are all
330 treated as 1--5 rating scales for this set of analyses. The null hypothesis of equal means for
331 the two groups would be true if respondents did not use the response labels. The second set
332 of analyses involves re-coding responses into matching categories, as shown in Table 4. Here
333 the null hypothesis corresponds to respondents not being influenced by whether the response
334 alternatives differentiate more at low or at high frequencies. This is important for
335 epidemiology because if these values differ it would produce different estimates for the
336 frequencies of these behaviors. The third set of analyses is exploratory. It relates to
337 potential moderators of these effects: self-reported political ideology and response time.
338 Whether which set of response alternative is presented affects responses to the attitude
339 questions is examined in the next section.

340 Figure 5 shows the histograms for the behavior questions when treated as 1--5 rating
341 scales. Table 5 shows the group means, the differences in means, effect sizes for these
342 differences (Cohen's d), and the 95% confidence intervals for these effect sizes. The values of

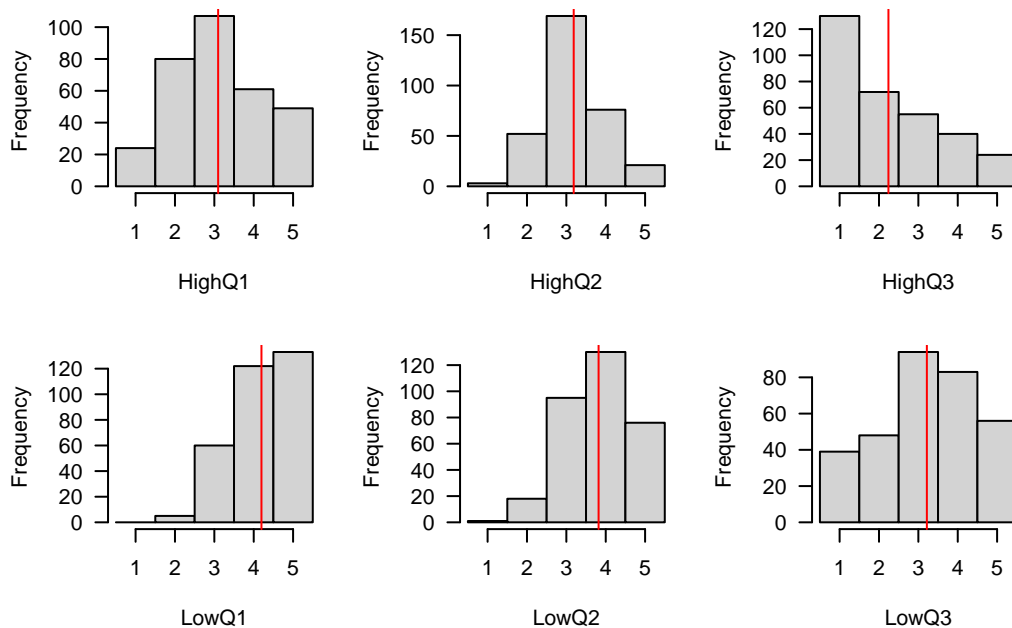


Figure 5

Histograms of the raw behavior responses, by condition. The red shows the condition mean for that variable.

343 Cohen's d are from 0.65 to 1.11. Cohen (1992) describes 0.5 as a *medium* effect and 0.8 as a
 344 *large* effect. Generic verbal labels for effect sizes can be problematic because the absolute
 345 meaning of any effect size is context dependent (Baguley, 2004; Lipsey et al., 2012). Here
 346 they are used to compare the relative size of these effects with those reported below when
 347 the response variables are treated as categorical.

348 Hotelling's T^2 (1931) is used to test if, as a group, the means of these three variables
 349 differed by condition. Box's M (1949), which tests equality of covariances, was statistically
 350 significant result, $\chi^2(6) = 60.38, p < .001$. Therefore a version of Hotelling's test that allows
 351 for heterogeneity of covariance matrices was used, and (as expected from the effect sizes) it
 352 was statistically significant: $F(3, 619.52) = 236.15, p < .001$.

353 The second set of analyses examines the re-coded responses within the dashed lines of
 354 Table 4. Table 6 shows the proportions for each of these categories for the two conditions.
 355 Cramér's V (with bias correction) is used here as the effect size measure. The BCa bootstrap
 356 intervals (2,000 replications) are shown. The χ^2 value is for testing the null hypothesis,

Table 5

Descriptive statistics for the raw values, 1--5, for the three behavior answers. This coding ignores the verbal labels, so these differences show some people pay attention to the labels.

	\bar{x}_{low}	\bar{x}_{high}	diff	LB	UB	Cohen's d	LB_d	UB_d
Hand wash frequency	4.20	3.10	-1.10	0.95	1.25	1.11	0.94	1.27
Hand wash duration	3.82	3.19	-0.63	0.50	0.76	0.75	0.59	0.91
Hand sanitizer frequency	3.22	2.24	-0.98	0.78	1.17	0.75	0.59	0.91

357 $V = 0$. Cohen (1988, p. 222) provides verbal labels for these effect sizes, and they vary by
 358 the degrees of freedom. For $df = 2$ (three categories), these are: small = .071, medium =
 359 .212, and large = .354. For $df = 3$ (four categories), these are: small = .058, medium = .173,
 360 and large = .289. Thus, as a group, using Cohen's terminology these show medium to large
 361 sized effects.

362 The third set of analyses concerns potential moderators. Two potential moderators of
 363 these effects are considered: self-reported political ideology and response time. Two politics
 364 variables were used. A political ideology variable is their raw score from 0--100. A political
 365 extremeness variable is created, which is the distance from the neutral response of 50, so a 0
 366 to 50 variable. The overall response time variable was skewed, 5.92. The natural logarithm
 367 was taken and the skewness was reduced to 0.66. t -tests were conducted on all three of these
 368 comparing the two conditions. None of these were statistically significant (unadjusted
 369 p -values shown): for political ideology: $t(632) = 1.63, p = .104$; for political extremeness:
 370 $t(632) = 1.81, p = .071$; and for logged response time: $t(639) = 0.76, p = .450$.

371 Moderation analyses were conducted on the behavior questions when treated as 1--5
 372 rating scales and the re-coded versions. MANOVAs were conducted predicting the three
 373 behavior questions as 1--5 rating scales. For each of the three moderator variables, the
 374 model with the interaction of the moderator with the condition (high versus low response
 375 alternative sets) was compared with the model with just the two main effects. The
 376 unadjusted p -values were: for political ideology: $p = .102$, for political extremeness: $p = .809$,
 377 and for response time: $p = .889$.

Table 6

Statistics comparing the recoded responses for the three behavior questions by condition.

Question	Cond.	Category				V (95% CI)	$\chi^2(df)$ (all $p < .001$)
		1	2	3	4		
Hand wash frequency	Low	.20	.38	.42		$V = .27$ (.19, .34)	$\chi^2(2) = 47.78$
	High	.07	.25	.68			
Hand wash duration	Low	.06	.30	.41	.24	$V = .22$ (.13, .28)	$\chi^2(3) = 31.85$
	High	.01	.16	.53	.30		
Hand sanitizer	Low	.57	.26	.17		$V = .22$ (.14, .29)	$\chi^2(2) = 31.82$
	High	.40	.22	.37			

378 The re-coded behavior questions are analyzed individually using multinomial logistic
 379 regression (using `multinom` from the **MASS** package from Venables and Ripley, 2002). The
 380 main effect of condition and moderator are included in a model, and this is compared with
 381 the model that also includes their interaction. With three behavior questions and three
 382 moderators, there are nine p -values. The unadjusted values ranged from .026 to .952. Only
 383 one of these is lower than the traditional $\alpha = .05$ level when not adjusted for multiple
 384 comparisons (those with extreme political views showed a smaller response alternative effect).
 385 Applying Holm's adjustment procedure for multiple comparisons, this becomes
 386 $p_{adjusted} = .238$. Following Spiegelhalter (2017), the results of exploratory analyses are
 387 reported but not elaborated upon.

388 *Attitude Questions*

389 The attitude questions were included to gauge whether influencing respondents to
 390 answer towards one end of the scales would influence how much they reported being worried
 391 about COVID-19 compared with others and how extreme their responses on the other
 392 attitude questions would be. As discussed in the introduction, the response alternatives may
 393 affect respondents about where they believe they place within the population norms. The
 394 skewness of the items ranged from -2.45 for seeking medical consultation to -0.54 for staying
 395 at home if not feeling well. The variables were transformed by ranking them and

Table 7*Pearson correlations among the norm-rank transformed attitude questions.*

	Econ	Health	Ucatch	FFcatch	Med	Home
Health	.595					
Ucatch	.357	.688				
FFcatch	.286	.590	.779			
Med	.047	.257	.255	.373		
Home	.370	.500	.419	.314	.234	
Science	.146	.374	.375	.402	.376	.319

396 normalizing the ranks (see van der Waerden, 1952; Wright, under review). The Pearson
 397 correlations of these transformed variables are shown in Table 7.

398 To increase statistical power, the pre-registered plan was if the items were correlated
 399 to combine them into a single dimension. The eigenvalues for the correlation matrix are:
 400 3.41, 1.13, 0.80, 0.62, 0.58, 0.28, and 0.19. Velicer, Eaton, and Fava (2000, see also
 401 Auerswald and Moshagen, 2019) discuss how to determine the number of dimensions. A
 402 common approach is the empirical Kaiser criterion, which compares the observed eigenvalues
 403 with those expected from random data (Braeken & van Assen, 2017). This suggests a single
 404 dimension. The first principal component was created and used for comparing groups. The
 405 t -test comparing means was non-significant: $t(639) = 0.87$, $p = .382$, Cohen's $d = 0.07$,
 406 95% $CI_d = (-0.09, 0.22)$.

407 Discussion and Recommendations

408 Evaluating public health campaigns and modeling disease spread requires estimating
 409 people's behaviors. This is usually done using surveys. Traditionally much effort has focused
 410 on the wording of question stems, and less on the response alternatives. Two studies focus
 411 on the response alternatives. The first shows that comparing responses across countries
 412 using vague quantifiers is difficult, but differences among the countries can almost
 413 completely be accounted for by the norms in that country. This provides strong support for
 414 the hypothesis put forth in Wright et al. (1994). While this stresses a difficulty making these

415 comparisons, because the country-differences can be accounted for by another variable, a
416 simple theory can be posited:

417 The meanings of the vague quantifiers are partially based on what respondents
418 believe the survey designers believe about the population norms. Much of the
419 respondents' beliefs will be based on the norms of their social and national
420 groups.

421 There will be specific linguistic nuances of particular words and geographic variation in their
422 usage. Further research would be necessary to construct improved ways, based on these
423 differences, to estimate numeric values. However, as noted in the recommendation section,
424 vague quantifiers should not be used if want to estimate or to compare groups with respect
425 to numeric values.

426 The second study was a randomized experiment to measure the effect of using
427 different response alternatives to estimate the frequency of health related behaviors. The
428 results showed that when several of the response alternatives were for high frequencies,
429 respondents gave answers corresponding to higher estimates than when they were for low
430 frequencies. The choice of response alternatives should be carefully considered when
431 estimating health related behaviors. Further, comparing estimates from surveys that use
432 different sets of response alternatives should be done with great caution, if at all. Which set
433 of response alternatives was presented did not have a statistically significant effect on the
434 attitudinal questions. This suggests that both sets were viewed as reasonable enough by
435 respondents not to create dissonance among them. This is taken into account in our
436 recommendations.

437 **Recommendations for Response Alternatives**

438 Behavioral frequency questions require respondents to answer questions about past
439 events. We make recommendations for two types of behaviors: rare and frequent. Whether
440 an event type is rare or frequent will depend on the sample and there will be overlap between
441 them, so survey designers should consider all the suggestions below where appropriate.

442 *Rare events*

443 Consider what is hopefully a rare event for the respondent, like a hospitalization,
444 catching COVID-19, or being laid off. It is likely for most respondents these are rare and it
445 is also likely that answers to these questions will be important both for the survey flow (e.g.,
446 on a COVID-19 survey if you answer YES to having COVID-19, you might be asked further
447 questions), and the estimates for these will be important for epidemiological models.
448 Therefore, accurate and precise answers are likely very important. The focus here is on the
449 response alternatives, but it is important to consider the memory limitations, even for rare
450 events (i.e., rare does not imply memorable), and that guiding the respondent using
451 procedures like the cognitive interview, as the term is used in eyewitness research, should be
452 considered (Fisher & Geiselman, 1992). In addition to remembering an event, respondents
453 usually need to provide information about *when* the event occurred. People have difficulty
454 saying when an event occurred (e.g., Friedman, 1993), and this is where issues about the
455 response format are most relevant.

456 Assuming that the respondents have time to respond to the survey and are near their
457 cell phone, respondents should be encouraged to consult resources (e.g., emails, texts,
458 vaccination card) to provide exact dates. If this information is not available, then
459 respondents could be presented with an *event history calendar*, as devised by Belli, Shay,
460 and Stafford (2001), and used in many health surveys (e.g., Schatz et al., 2020). These allow
461 respondents to fill in notable events, like a child's graduation and holidays, onto the
462 calendars, and then allow the respondent to think where the event in question happened in
463 relation to these. If using the event history calendar is impractical, most online surveys allow
464 a calendar response so the respondent can give a precise date, and this can also be used if
465 there were multiple incidences in the time frame. For long duration events a start and end
466 could be provided, and respondents could also provide a range if they are uncertain.

Frequent Events

By frequent events we mean those that, for most respondents, are likely to occur multiple times each week (e.g., handwashing, eating a piece of fruit). As with rare events it is important to consider the cognitive limitations of the respondent, and in particular whether the respondent is likely to use some estimation heuristic or try to recall and count all episodes. Each of these has potential biases (e.g., people are more likely to not remember an event than to create a false memory for a non-existent event, e.g., Wright, Loftus, and Hall, 2001, so recalling each incident is likely to result in an under-estimate), so if the accuracy is of much importance diary methods and experience-sampling methods (e.g., K. Xie, Heddy, & Vongkulluksn, 2019) can be used, though these require coordination with respondents prior to the survey and in the case of experience-sampling methods more technology. For some frequent events, there might be available resources available retrospectively (e.g., filtering through the trash for some food consumption events), though these would not be available for most event types and would be more effort than is likely appropriate. The following assumes these resources are not available and that this is for a single retrospective survey.

Survey designers can be interested in both well-defined events and those that are not well-defined. While standard practice encourages survey questions to be well-defined, if the goal is to compare groups on, for example, how worried they are about catching COVID-19, or some other psychological state, there is often no way to make these precise. Because of this, it will be difficult to interpret numeric estimates of the frequencies with much confidence unless it can be made clear to respondents what, in this case, an episode of worry would be. In these situations it can be prudent to use vague quantifiers. The vagueness of the response alternative matches with the event. Often non-numeric responses can also be used. For example, comparison questions can also be used when this matches what people's beliefs about the events are like. This would be a situation where pre-testing using think-aloud protocols would likely be valuable to show how respondents think about these event descriptions (Willis, 2005, especially Chapter 4).

494 For well-defined events, vague quantifiers should be avoided. The choice should be
495 between free recall procedures and a set of numeric response alternatives. We qualify free
496 recall with *ish* to emphasize that this may not mean just writing something in a box. For
497 online surveys, the software can force respondents to enter text in a specified format (e.g., as
498 a number if that is necessary for subsequent analysis, rather than accepting, for example, a
499 respondent writing “about 4 to 8, maybe”). An issue with this is that people’s memories
500 may not be that precise. Goldsmith, Koriat, and Weinberg-Eliezer (2002) describe how
501 people’s memory granularity varies for different memories. This can be accounted for when
502 you ask for numeric information by allowing people to provide a range of values
503 corresponding to their confidence (e.g., Weber & Brewer, 2008). Alternatively, sets of
504 numeric response alternatives can be used, though this has the limitation that the
505 alternatives will usually be a range of values so may not be as precise as free recall methods.
506 However, study 2 showed the the choice of response alternatives can make a difference. We
507 recommend using a large number of alternatives that account for how different groups
508 within the sample will have different expectations.

509 **Summary**

510 Survey data are used in social science research for informing economists about
511 consumer behaviors, for health researchers evaluating compliance with different campaigns,
512 for sociologists and psychologists constructing theories of why people behave in they ways
513 that they do, and for many other purposes. Behavioral frequency questions have a special
514 place within survey methods as researchers often translate responses into numeric estimates
515 for the behaviors, and sometimes the precision of these estimates is critical (e.g., to
516 epidemiologists predicting trends for the COVID-19 pandemic). While people often focus on
517 the way the event itself is described in the question stem, there is less focus on the response
518 alternatives. We focus on the response alternatives.

519 Our first study examined how people, across thirty countries, answered questions
520 about hand washing and hand sanitizing. We found that people in different countries

521 interpreted the vague quantifiers used as response alternatives differently. We were able to
522 account for most of the differences among countries using estimates of the behavior in the
523 countries from a different set of respondents. We do not recommend using vague quantifiers
524 with relatively well defined events like hand washing, but allowing people to provide
525 numerical estimates. Our second study showed that care is still necessary when doing this.
526 Using different sets of response alternatives produced different estimates of the behavior.
527 One consequence of this is that comparisons between studies that use different sets of
528 response alternatives should be done cautiously, if at all.

529 One conclusion from our studies is that the choice of response alternatives should be
530 carefully considered and the deciding how to construct them may be difficult. It may require
531 careful pilot research and techniques like cognitive interviewing and in particular think-aloud
532 protocols (Willis, 2005). We provide a list of recommendations to allow researchers to start
533 thinking about their choices of response alternatives for behavioral frequency questions.

References

- 534
- 535 Auerswald, M., & Moshagen, M. (2019). How to determine the number of factors to retain
536 in exploratory factor analysis: A comparison of extraction methods under realistic
537 conditions. *Psychological Methods, 24*, 468--491. doi: 10.1037/met0000200
- 538 Baguley, T. S. (2004). Understanding statistical power in the context of applied research.
539 *Applied Ergonomics, 35*, 73--80. doi: 10.1016/j.apergo.2004.01.002
- 540 Belli, R. F., Conrad, F. G., & Wright, D. B. (2007). Cognitive psychology and survey
541 methodology: nurturing the continuing dialogue between disciplines. *Applied Cognitive
542 Psychology, 21*, 141--144. doi: 10.1002/acp.1333
- 543 Belli, R. F., Shay, W. L., & Stafford, F. P. (2001). Event history calendars and question list
544 surveys: A direct comparison of interviewing methods. *Public Opinion Quarterly, 65*,
545 45--74.
- 546 Box, G. E. P. (1949). A general distribution theory for a class of likelihood criteria.
547 *Biometrika, 36*, 317--346. doi: 10.1093/biomet/36.3-4.317
- 548 Braeken, J., & van Assen, M. A. L. M. (2017). An empirical Kaiser criterion. *Psychological
549 Methods, 22*, 450 --466. doi: 10.1037/met0000074
- 550 Burton, S., & Blair, E. (1991). Task conditions, response formulation processes, and
551 response accuracy for behavioral frequency questions in surveys. *Public Opinion
552 Quarterly, 55*, 50--79. doi: 10.1086/269241
- 553 Cantril, H. (1965). *The pattern of human concerns*. New Brunswick, NJ: Rutgers University
554 Press.
- 555 Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Hillsdale,
556 NJ: Lawrence Erlbaum Associates.
- 557 Cohen, J. (1992). A power primer. *Psychological Bulletin, 112*, 155--159. doi:
558 10.1037/0033-2909.112.1.155

- 559 Fienberg, S. E., Loftus, E. F., & Tanur, J. M. (1985). Cognitive aspects of health survey
560 methodology: An overview. *The Milbank Memorial Fund Quarterly. Health and*
561 *Society, 63*, 547--564.
- 562 Fisher, R. P., & Geiselman, R. E. (1992). *Memory-enhancing techniques for investigative*
563 *interviewing: The cognitive interview*. Springfield, IL: Charles C. Thomas, Publisher.
- 564 Friedman, W. J. (1993). Memory for the time of past events. *Psychological Bulletin, 113*,
565 44--66. doi: 10.1037/0033-2909.113.1.44
- 566 Gadarian, S. K., Goodman, S. W., & Pepinsky, T. B. (2021). Partisanship, health behavior,
567 and policy attitudes in the early stages of the COVID-19 pandemic. *PLoS ONE, 16*(4),
568 e0249596. doi: 10.1371/journal.pone.0249596
- 569 Gaskell, G. D., O'Muircheartaigh, C. A., & Wright, D. B. (1994). Survey questions about
570 the frequency of vaguely defined events: the effects of response alternatives. *Public*
571 *Opinion Quarterly, 58*, 241--254. doi: 10.1086/269420
- 572 Goldsmith, M., Koriat, A., & Weinberg-Eliezer, A. (2002). Strategic regulation of grain size
573 memory reporting. *Journal of Experimental Psychology: General, 131*, 73--95.
- 574 Grice, P. (1975). Logic and conversation. In P. Cole & J. Morgan (Eds.), *Syntax and*
575 *semantics. 3: Speech acts* (pp. 41--58). New York: Academic Press.
- 576 Groves, R. M., Fowler, F. J., Jr., Couper, M. P., Lepkowski, J. M., Singer, E., &
577 Tourangeau, R. (2009). *Survey methodology* (2nd ed.). Hobokon, NJ: Wiley.
- 578 Hotelling, H. (1931). The generalization of Student's ratio. *Annals of Mathematical*
579 *Statistics, 2*, 360--378. doi: 10.1214/aoms/1177732979
- 580 Huttenlocher, J., Hedges, L. V., & Duncan, S. (1991). Categories and particulars: Prototype
581 effects in estimating spatial location. *Psychological Review, 98*, 352--376. doi:
582 10.1037/0033-295X.98.3.352
- 583 Jones, S. P. (2020). *Imperial College London YouGov Covid data hub, v1.0*. Imperial
584 College London Big Data Analytical Unit and YouGov Plc.

- 585 Lee, L. M., & Thacker, S. B. (2011). Public health surveillance and knowing about health in
586 the context of growing sources of health data. *American Journal of Preventive*
587 *Medicine*, *41*, 636--640. doi: 10.1016/j.amepre.2011.08.015
- 588 Lipsey, M. W., Puzio, K., Yun, C., Hebert, M. A., Steinka-Fry, K., Cole, M. W., . . . Busick,
589 M. D. (2012). *Translating the statistical representation of the effects of education*
590 *interventions into more readily interpretable forms. (NCSEER 2013-3000)*. Washington,
591 DC: National Center for Special Education Research, Institute of Education Sciences,
592 U.S. Department of Education. Retrieved from
593 <https://ies.ed.gov/ncser/pubs/20133000/pdf/20133000.pdf>
- 594 Loftus, E. F., Fienberg, S. E., & Tanur, J. M. (1985). Cognitive psychology meets the
595 national survey. *American Psychologist*, *40*, 175--180. doi:
596 10.1037/0003-066X.40.2.175
- 597 Loftus, E. F., & Marburger, W. (1983). Since the eruption of Mt. St. Helens, has anyone
598 beaten you up? improving the accuracy of retrospective reports with landmark events.
599 *Memory & Cognition*, *11*, 114--120. doi: 10.3758/BF03213465
- 600 Neter, J., & Waksberg, J. (1964). A study of response errors in expenditures data from
601 household interviews. *Journal of the American Statistical Association*, *59*, 18--55. doi:
602 10.1080/01621459.1964.10480699
- 603 Schaeffer, N. C. (1991). Hardly ever or constantly? Group comparisons using vague
604 quantifiers. *Public Opinion Quarterly*, *55*, 395--423. doi: 10.1086/269270
- 605 Schaeffer, N. C., & Dykema, J. (2020). Advances in the science of asking questions. *Annual*
606 *Review of Sociology*, *46*, 37--60. doi: 10.1146/annurev-soc-121919-054544
- 607 Schatz, E., Knight, L., Belli, R. F., & Mojola, S. A. (2020). Assessing the feasibility of a life
608 history calendar to measure HIV risk and health in older South Africans. *PLoS: One*,
609 *15*(1), e0226024. doi: 10.1371/journal.pone.0226024

- 610 Schwarz, N. (1995). What respondents learn from questionnaires: The survey interview and
611 the logic of conversation. *International Statistical Review / Revue Internationale de*
612 *Statistique*, 63, 153--168.
- 613 Schwarz, N. (2010). Measurement as cooperative communication: What research
614 participants learn from questionnaires. In G. Walford, M. Viswanathan, & E. Tucker
615 (Eds.), *The SAGE handbook of measurement* (pp. 43--59). Thousand Oaks, CA: Sage.
- 616 Schwarz, N., & Hippler, H.-J. (1987). What response scales may tell your respondents:
617 Informative functions of response alternatives. In H.-J. Hippler, N. Schwarz, &
618 S. Sudman (Eds.), *Social information processing and survey methodology* (pp.
619 163--178). New York, NY: Springer.
- 620 Schwarz, N., Hippler, H.-J., Deutsch, B., & Strack, F. (1985). Response scales: effects of
621 category range on reported behavior and comparative judgments. *Public Opinion*
622 *Quarterly*, 49, 388--395. doi: 10.1086/268936
- 623 Schwarz, N., Strack, F., Müller, G., & Chassein, B. (1988). The range of response
624 alternatives may determine the meaning of the question: further evidence on
625 informative functions of response alternatives. *Social Cognition*, 6, 107-117. doi:
626 10.1521/soco.1988.6.2.107
- 627 Spiegelhalter, D. J. (2017). Trust in numbers. *Journal of the Royal Statistical Society:*
628 *Series A (Statistics in Society)*, 180, 948--965. doi: 10.1111/rssa.12302
- 629 Sudman, S., & Bradburn, N. M. (1973). Effects of time and memory factors on response in
630 surveys. *Journal of the American Statistical Association*, 68, 805--815.
- 631 Thompson, C. P., Skowronski, J. J., & Lee, D. J. (1988). Telescoping in dating naturally
632 occurring events. *Memory & Cognition*, 16, 461--468. doi: 10.3758/BF03214227
- 633 Tourangeau, R. (2003). Cognitive aspects of survey measurement and mismeasurement.
634 *International Journal of Public Opinion Research*, 15, 3--7. doi: 10.1093/ijpor/15.1.3

- 635 US Census Bureau. (2021). *Household pulse survey: measuring social and economic impacts*
636 *during the Coronavirus pandemic*. Retrieved from
637 www.census.gov/programs-surveys/household-pulse-survey.html
- 638 van der Waerden, B. L. (1952). Order tests for the two-sample problem and their power.
639 *Indagationes Mathematicae (Proceedings)*, 55, 453--458. doi:
640 10.1016/S1385-7258(52)50063-5
- 641 Velicer, W. F., Eaton, C. A., & Fava, J. L. (2000). Construct explication through factor or
642 component analysis: a review and evaluation of alternative procedures for determining
643 the number of factors or components. In R. D. Goffin & E. Helmes (Eds.), *Problems*
644 *and solutions in human assessment* (pp. 41--71). Boston, MA: Kluwer.
- 645 Venables, W. N., & Ripley, B. D. (2002). *Modern applied statistics with S* (Fourth ed.).
646 New York: Springer. Retrieved from <http://www.stats.ox.ac.uk/pub/MASS4>
647 (ISBN 0-387-95457-0)
- 648 Weber, N., & Brewer, N. (2008). Eyewitness recall: Regulation of grain size and the role of
649 confidence. *Journal of Experimental Psychology: Applied*, 14, 50--60. doi:
650 10.1037/1076-898X.14.1.50
- 651 Willis, G. B. (2005). *Cognitive interviewing: tools for improving questionnaire design*. Sage
652 Publications, Inc. doi: 10.4135/9781412983655
- 653 Wilson, D., & Sperber, D. (2012). *Meaning and relevance*. Cambridge University Press.
- 654 Wright, D. B. (under review). A robust alternative to Pearson's correlation for testing
655 associations and use in latent variable models.
- 656 Wright, D. B., Gaskell, G. D., & O'Muircheartaigh, C. A. (1994). How much is 'quite a bit'?
657 Mapping between numerical values and vague quantifiers. *Applied Cognitive*
658 *Psychology*, 8, 479--496. doi: <https://doi.org/10.1002/acp.2350080506>
- 659 Wright, D. B., Gaskell, G. D., & O'Muircheartaigh, C. A. (1997). How response alternatives
660 affect different kinds of behavioural frequency questions. *British Journal of Social*
661 *Psychology*, 36, 443--456. doi: 10.1111/j.2044-8309.1997.tb01143.x

- 662 Wright, D. B., Loftus, E. F., & Hall, M. (2001). Now you see it; now you don't: inhibiting
663 recall and recognition of scenes. *Applied Cognitive Psychology*, *15*, 471--482. doi:
664 10.1002/acp.719
- 665 Xie, K., Heddy, B. C., & Vongkulluksn, V. W. (2019). Examining engagement in context
666 using experience-sampling method with mobile technology. *Contemporary Educational*
667 *Psychology*, *59*, 101788. doi: <https://doi.org/10.1016/j.cedpsych.2019.101788>
- 668 Xie, Y. (2013). *Dynamic documents with R and **knitr***. Boca Raton, FL: Chapman and
669 Hall/CRC.