



Kao, A. B., Berdahl, A. M., Hartnett, A. T., Lutz, M. J., Bak-Coleman, J. B., Ioannou, C. C., ... Couzin, I. D. (2018). Counteracting estimation bias and social influence to improve the wisdom of crowds. *Journal of the Royal Society Interface*, 15(141), [20180130].  
<https://doi.org/10.1098/rsif.2018.0130>

Peer reviewed version

Link to published version (if available):  
[10.1098/rsif.2018.0130](https://doi.org/10.1098/rsif.2018.0130)

[Link to publication record in Explore Bristol Research](#)  
PDF-document

This is the accepted author manuscript (AAM). The final published version (version of record) is available online via The Royal Society at DOI: 10.1098/rsif.2018.0130. Please refer to any applicable terms of use of the publisher.

## University of Bristol - Explore Bristol Research

### General rights

This document is made available in accordance with publisher policies. Please cite only the published version using the reference above. Full terms of use are available:  
<http://www.bristol.ac.uk/pure/about/ebr-terms>

# Counteracting estimation bias and social influence to improve the wisdom of crowds

Albert B. Kao<sup>a,\*</sup>, Andrew Berdahl<sup>b</sup>, Andrew T. Hartnett<sup>c</sup>, Matthew J. Lutz<sup>d</sup>, Joseph Bak-Coleman<sup>e</sup>, Christos C. Ioannou<sup>f</sup>, Xingli Giam<sup>g</sup>, Iain D. Couzin<sup>d,h</sup>

<sup>a</sup>*Department of Organismic and Evolutionary Biology, Harvard University, Cambridge, MA, USA*

<sup>b</sup>*Santa Fe Institute, Santa Fe, NM, USA*

<sup>c</sup>*Argo AI, Pittsburgh, PA, USA*

<sup>d</sup>*Department of Collective Behaviour, Max Planck Institute for Ornithology, Konstanz, Germany*

<sup>e</sup>*Department of Ecology and Evolutionary Biology, Princeton University, Princeton, NJ, USA*

<sup>f</sup>*School of Biological Sciences, University of Bristol, Bristol, UK*

<sup>g</sup>*Department of Ecology and Evolutionary Biology, University of Tennessee, Knoxville, TN, USA*

<sup>h</sup>*Chair of Biodiversity and Collective Behaviour, Department of Biology, University of Konstanz, Konstanz, Germany*

---

## Abstract

Aggregating multiple non-expert opinions into a collective estimate can improve accuracy across many contexts. However, two sources of error can diminish collective wisdom: individual estimation biases and information sharing between individuals. Here we measure individual biases and social influence rules in multiple experiments involving hundreds of individuals performing a classic numerosity estimation task. We first investigate how existing aggregation methods, such as calculating the arithmetic mean or the median, are influenced by these sources of error. We show that the mean tends to overestimate, and the median underestimate, the true value for a wide range of numerosities. Quantifying estimation bias, and mapping individual bias to collective bias, allows us to develop and validate three new aggregation measures that effectively counter sources of collective estimation error. In addition, we present results from a further experiment that quantifies the social influence rules that individuals employ when incorporating personal estimates with social information. **[We show that the corrected mean is remarkably robust to social influence, retaining high accuracy in the presence or absence of social influence, across numerosities, and across different methods for averaging**

---

\*Author for correspondence: Albert B. Kao, e-mail: albert.kao@gmail.com

**social information.]** Utilizing knowledge of estimation biases and social influence rules may therefore be an inexpensive and general strategy to improve the wisdom of crowds.

*Keywords:* wisdom of crowds, collective intelligence, social influence, estimation bias, numerosity

---

## 1 1. Introduction

2 The proliferation of online social platforms has enabled the rapid expression of opinions  
3 on topics as diverse as the outcome of political elections, policy decisions, or the future  
4 performance of financial markets. Because non-experts contribute the majority of these  
5 opinions, they may be expected to have low predictive power. However, it has been shown  
6 empirically that by aggregating these non-expert opinions, usually by taking the arith-  
7 metic mean or the median of the set of estimates, the resulting ‘collective’ estimate can be  
8 highly accurate [1–6]. Experiments with non-human animals have demonstrated similar  
9 results [7–12], suggesting that aggregating diverse estimates can be a simple strategy for  
10 improving estimation accuracy across contexts and even species.

11 Theoretical explanations for this ‘wisdom of crowds’ typically invoke the law of large  
12 numbers [1, 13, 14]. If individual estimation errors are unbiased and center at the true  
13 value, then averaging the estimates of many individuals will increasingly converge on  
14 the true value. However, empirical studies of individual human decision-making readily  
15 contradict this theoretical assumption. A wide variety of cognitive and perceptual biases  
16 have been documented in which humans seemingly deviate from rational behavior [15–17].  
17 Empirical ‘laws’ such as the Stevens’ power law [18] have described the non-linear rela-  
18 tionship between the actual magnitude, and subjective perception, of a physical stimulus.  
19 Such nonlinearities can lead to a systematic under- or over-estimation of a stimulus, as  
20 is frequently observed in numerosity estimation tasks [19–22]. Furthermore, the Weber-  
21 Fechner law [23] implies that log-normal, rather than normal, distributions of estimates  
22 are common. When such biased individual estimates are aggregated, the resulting collec-  
23 tive estimate may also be biased, although the mapping between individual and collective  
24 biases is not well understood.

25 Sir Francis Galton was one of the first to consider the effect of biased opinions on  
26 the accuracy of collective estimates. He preferred the median over the arithmetic mean,  
27 arguing that the latter measure “give[s] a voting power to ‘cranks’ in proportion to their  
28 crankiness” [24]. However, if individuals are prone to under- or over-estimation in a par-  
29 ticular task, then the median will also under- or over-estimate the true value. Other ag-  
30 gregation measures have been proposed to improve the accuracy of the collective estimate,  
31 such as the geometric mean [25], the ‘trimmed mean’ (where the tails of a distribution  
32 of estimates are trimmed and then the arithmetic mean is calculated from the resulting  
33 truncated distribution) [26], and the average of the arithmetic mean and median [27].  
34 Although these measures may empirically improve accuracy in some cases, they tend not  
35 to address directly the root causes of collective error (*i.e.*, estimation bias). Therefore, it  
36 is not well understood how they generalize to other contexts and how close they are to  
37 the optimal aggregation strategy.

38 Many (though not all) models of the wisdom of crowds also assume that opinions  
39 are generated independently of one another, which tends to maximize the information  
40 contained within the set of opinions [1, 13, 14]. But in real world contexts, it is more  
41 common for individuals to share information with, and influence, one another [25, 28]. In  
42 such cases, the individual estimates used to calculate a collective estimate will be corre-  
43 lated to some degree. Social influence can not only shrink the distribution of estimates [25]  
44 but may also systematically shift the distribution, depending on the rules that individuals  
45 follow when updating their personal estimate in response to available social information.  
46 For example, if individuals with extreme opinions are more resistant to social influence,  
47 then the distribution of estimates will tend to shift towards these opinions, leading to  
48 changes in the collective estimate as individuals share information with each other. In  
49 short, social influence may induce estimation bias, even if individuals in isolation are  
50 unbiased.

51 Quantifying how both individual estimation biases and social influence affect collective  
52 estimation is therefore crucial to optimizing, and understanding the limits of, the wisdom  
53 of crowds. Such an understanding would help to identify which of the existing aggregation

54 measures should lead to the highest accuracy. It could also permit the design of novel  
55 aggregation measures that counteract these major sources of error, potentially improving  
56 both the accuracy and robustness of the wisdom of crowds beyond that allowed by existing  
57 measures.

58 Here, we collected five new datasets, and analyzed eight existing datasets from the lit-  
59 erature, to characterize individual estimation bias in a well-known wisdom of crowds task,  
60 the ‘jellybean jar’ estimation problem. In this task, individuals in isolation simply esti-  
61 mate the number of objects (such as jellybeans, gumballs, or beads) in a jar [5, 6, 29, 30]  
62 (see Methods for details). We then performed an experiment manipulating social infor-  
63 mation to quantify the social influence rules that individuals use during this estimation  
64 task (Methods). We used these results to quantify the accuracy of a variety of aggregation  
65 measures, and identified new aggregation measures to improve collective accuracy in the  
66 presence of individual bias and social influence.

## 67 **2. Methods**

### 68 *2.1. Numerosity estimation*

69 For the five datasets that we collected, we recruited members of the community in  
70 Princeton, NJ, USA on April 26–28 and May 1, 2012, and in Santa Fe, NM, USA on  
71 October 17–20, 2016. Each participant was presented with one jar containing one of  
72 the following numbers of objects: 54 ( $n = 36$ ), 139 ( $n = 51$ ), 659 ( $n = 602$ ), 5897  
73 ( $n = 69$ ), or 27852 ( $n = 54$ ) (see Figure 1a for a representative photograph of the kind of  
74 object and jar used for the three smallest numerosities, and Figure S1 for a representative  
75 photograph of the kind of object and jar used for the largest two numerosities.). To  
76 motivate accurate estimates, the participants were informed that the estimate closest to  
77 the true value for each jar would earn a monetary prize. The participants then estimated  
78 the number of objects in the jar. No time limit was set, and participants were advised  
79 not to communicate with each other after completing the task.

80 Eight additional datasets were included for comparative purposes and were obtained  
81 from refs. [5, 6, 29, 30].

82 Details of statistical analyses and simulations performed on these data are provided  
83 in the electronic supplementary material.

### 84 *2.2. Social influence experiment*

85 For the experiments run in Princeton (number of objects  $J = 659$ ), we additionally  
86 tested the social influence rules that individuals use. The participants first recorded their  
87 initial estimate,  $G_1$ . Next, participants were given ‘social’ information, in which they  
88 were told that  $N = \{1, 2, 5, 10, 50, 100\}$  previous participants’ estimates were randomly  
89 selected and that the ‘average’ of these guesses,  $S$ , was displayed on a computer screen.  
90 Unbeknownst to the participant, this social information was artificially generated by the  
91 computer, allowing us to control, and thus decouple, the perceived social group size and  
92 social distance relative to the participant’s initial guess. Half of the participants were  
93 randomly assigned to receive social information taken from a uniform distribution from  
94  $G_1/2$  to  $G_1$ , and the other half received social information from a uniform distribution  
95 from  $G_1$  to  $2G_1$ . Participants were then given the option to revise their initial guess by  
96 making a second estimate,  $G_2$ , based on their personal estimate and the perceived social  
97 information that they were given. Participants were informed that only the second guess  
98 would count toward winning a monetary prize. We therefore controlled the social group  
99 size by varying  $N$  and controlled the social distance independently of the participant’s  
100 accuracy by choosing  $S$  from  $G_1/2$  to  $2G_1$ .

101 Details of the social influence model and simulations performed on these data are  
102 provided in the electronic supplementary material.

### 103 *2.3. Designing ‘corrected’ aggregation measures*

104 For a log-normal distribution, the expected value of the mean is given by  $X_{\text{mean}} =$   
105  $\exp(\mu + \sigma^2/2)$  and the expected value of the median is  $X_{\text{median}} = \exp(\mu)$ , where  $\mu$  and  
106  $\sigma$  are the two parameters describing the distribution. Our empirical measurements of esti-  
107 mation bias resulted in the best-fit relationships  $\mu = m_\mu \ln(J) + b_\mu$  and  $\sigma = m_\sigma \ln(J) + b_\sigma$   
108 (Figure 1c-d). We replace  $\mu$  and  $\sigma$  in the first two equations with the best-fit relation-  
109 ships, and then solve for  $J$ , which is our new, ‘corrected’, estimate of the true value. This

110 results in a ‘corrected’ arithmetic mean:

$$X_{\text{mean}}^C = \exp \left( \left( \sqrt{2m_\sigma^2 (\ln X_{\text{mean}} - b_\mu) + 2m_\mu^2 \left( \frac{1}{2} + \frac{m_\sigma b_\sigma}{m_\mu} \right)} - (m_\sigma b_\sigma + m_\mu) \right) / m_\sigma^2 \right)$$

111 and a ‘corrected’ median:

$$X_{\text{median}}^C = \exp ((\ln X_{\text{median}} - b_\mu) / m_\mu)$$

112 This procedure can be readily adapted for other estimation tasks, distributions of  
113 estimates, and estimation biases.

#### 114 *2.4. A maximum-likelihood aggregation measure*

115 For this aggregation measure, the full set of estimates is used to form a new collective  
116 estimate, rather than using just the mean or the median to generate a corrected measure.  
117 We again invoke the best-fit relationships in Figure 1c-d, which imply that, for a given  
118 actual number of objects  $J$ , we expect a log-normal distribution described by parameters  
119  $\mu = m_\mu \ln(J) + b_\mu$  and  $\sigma = m_\sigma \ln(J) + b_\sigma$ . We therefore scan across values of  $J$  and  
120 calculate the likelihood that each associated log-normal distribution generated the given  
121 set of estimates. The numerosity that maximizes this likelihood is the collective estimate  
122 of the true value.

### 123 **3. Results**

#### 124 *3.1. Quantifying estimation bias*

125 To uncover individual biases in estimation tasks, we first sought to characterize how  
126 the distribution of individual estimates changes as a function of the true number of objects  
127  $J$  (Figure 1a). We performed experiments across a >500-fold range of numerosities, from  
128 54 to 27852 objects, with a total of 812 people sampled across the experiments. For all  
129 numerosities tested, an approximately log-normal distribution was observed (see Figure  
130 1b for a histogram of an example dataset, Figure S2 for histograms of all other datasets,  
131 and Figure S3 for a comparison of the datasets to log-normal distributions). Log-normal  
132 distributions can be described by two parameters,  $\mu$  and  $\sigma$ , which correspond to the

133 arithmetic mean and standard deviation, respectively, of the normal distribution that  
134 results when the original estimates are log-transformed (Figure 1b, inset, and section 1  
135 of the electronic supplementary material on how the maximum-likelihood estimates of  $\mu$   
136 and  $\sigma$  were computed for each dataset).

137 We found that the shape of the log-normal distribution changes in a predictable man-  
138 ner as the numerosity changes. In particular, the two parameters of the log-normal dis-  
139 tribution,  $\mu$  and  $\sigma$ , both exhibit a linear relationship with the logarithm of the number  
140 of objects in the jar (Figure 1c-d). These relationships hold across the entire range of  
141 numerosities that we tested (which spans nearly three orders of magnitude). That the  
142 parameters of the distribution co-vary closely with numerosity allows us to directly com-  
143 pute how the magnitude of various aggregation measures changes with numerosity, and  
144 provides us with information about human estimation behavior which we can exploit to  
145 improve the accuracy of the aggregation measures.

### 146 *3.2. Expected error of aggregation measures*

147 We used the maximum-likelihood relationships shown in Figure 1c-d to first compute  
148 the expected value of the arithmetic mean, given by  $\exp(\mu + \sigma^2/2)$ , and the median, given  
149 by  $\exp(\mu)$ , of the log-normal distribution of estimates, across the range of numerosities  
150 that we tested empirically (between 54 and 27852 objects). We then compared the magni-  
151 tude of these two aggregation measures to the true value to identify any systematic biases  
152 in these measures (we note that any aggregation measure may be examined in this way,  
153 but for clarity here we display just the two most commonly used measures).

154 Overall, across the range of numerosities tested, we found that the arithmetic mean  
155 tended to overestimate, while the median tended to underestimate, the true value (Figure  
156 2a). This is corroborated by our empirical data: for four out of the five datasets, the  
157 mean overestimated the true value, while the median underestimated the true value in  
158 four of five datasets (Figure 2a). **[We note that our model predicts qualitatively**  
159 **different patterns for very small numerosities (outside of the range that we**  
160 **tested experimentally). Specifically, in this regime the model predicts that**  
161 **the mean and the median both overestimate the true value, with large relative**



162 errors for both measures. However, we expect humans to behave differently  
163 when presented with a small number of objects that can be counted directly  
164 compared to uncountably many objects; therefore, we avoid extrapolating our  
165 results and apply our model only on the range that we tested experimentally  
166 (spanning nearly three orders of magnitude).]

167 That the median tends to underestimate the true value implies that the majority of  
168 individuals underestimate the true numerosity. This conforms with the results of other  
169 studies demonstrating an underestimation bias in numerosity estimation in humans (*e.g.*,  
170 [20–22, 31]). Despite this, the arithmetic mean tends to overestimate the true value be-  
171 cause the log-normal distribution has a long tail (Figure 1b), which inflates the mean.  
172 Indeed, because the parameter  $\sigma$  increases with numerosity, the dispersion of the distri-  
173 bution is expected to increase disproportionately quickly with numerosity, such that the  
174 coefficient of variation (the ratio between the standard deviation and the mean of the un-  
175 transformed estimates) increases with numerosity (Figure S4). This finding differs from  
176 other results showing a constant coefficient of variation across numerosities [19, 20]. This  
177 contrasting result may be explained by the larger-than-typical range of numerosities that  
178 we evaluated here (with respect to previous studies), which improves our ability to detect  
179 a trend in the coefficient of variation. Alternatively (and not mutually exclusively), it may  
180 result from other studies displaying many numerosities to the same participant, which may  
181 cause correlations in a participant’s estimates [20, 21] and reduce variation. By contrast,  
182 we only showed a single jar to each participant in our estimation experiments. Overall,  
183 the degree of underestimation and overestimation of the median and mean, respectively,  
184 was approximately equal across the range of numerosities tested, and we did not detect  
185 consistent differences in accuracy between these two aggregation measures (Figure 2b).

### 186 3.3. Designing and testing aggregation measures that counteract estimation bias

187 Knowing the expected error of the aggregation measures relative to the true value,  
188 we can design new measures to counter this source of collective estimation error. Using  
189 this methodology, we specify functional forms of the ‘corrected’ arithmetic mean and the  
190 ‘corrected’ median (Methods). In addition to these two adjusted measures, we propose a

191 maximum-likelihood method that uses the full set of estimates, rather than just the mean  
192 or median, to locate the numerosity that is most likely to have produced those estimates  
193 (Methods). Although applied here to the case of log-normal distributions and particular  
194 relationships between numerosity and the parameters of the distributions, our procedure  
195 is general and could be used to construct specific corrected measures appropriate for other  
196 distributions and relationships, subsequent to empirically characterizing these patterns.

197     Once the corrected measures have been parameterized for a specific context, they can  
198 be applied to a new test dataset to produce a improved collective estimate from that  
199 data. However, the three new measures are predicted to have near-zero error only in their  
200 expected values, which assumes an infinitely large test dataset (and that the corrected  
201 measures have been accurately parameterized). A finite-sized set of estimates, on the other  
202 hand, will generally exhibit some deviation from the expected value. It is possible that the  
203 measures will produce different noise distributions around the expected value, which will  
204 affect their real-world accuracy. To address this, we measured the overall accuracy of the  
205 aggregation measures across a wide range of test sample sizes and numerosities, simulating  
206 datasets by drawing samples using the maximum-likelihood fits shown in Figure 1c-d. We  
207 also conducted a separate analysis, in which we generate test datasets by drawing samples  
208 directly from our experimental data, the results of which we include in the electronic  
209 supplementary material (see section 2 of the electronic supplementary material for details  
210 on both methodologies and for justification of why we chose to include the results from  
211 the simulated data in the main text.)

212     We compared each of the new aggregation measures to the arithmetic mean, the  
213 median, and three other ‘standard’ measures that have been described previously in the  
214 literature: the geometric mean, the average of the mean and the median, and a trimmed  
215 mean (where we remove the smallest 10% of the data, and the largest 10% of the data,  
216 before computing the arithmetic mean), in pairwise fashion, calculating the fraction of  
217 simulations in which one measure had lower error than the other.

218     All three new aggregation measures outperformed all of the other measures (Figure  
219 3a, left five columns), displaying lower error in 58–78% of simulations. Comparing the

220 three new measures against each other, the maximum-likelihood measure performed best,  
221 followed by the corrected mean, while the corrected median resulted in the lowest overall  
222 accuracy (Figure 3a, right three columns). The 95% confidence intervals of the percentages  
223 are, at most,  $\pm 1\%$  of the stated percentages (binomial test,  $n = 10000$ ), and therefore the  
224 results shown in Figure 3a are all significantly different from chance. The results from our  
225 alternate analysis, using samples drawn from our experimental data, are broadly similar,  
226 albeit somewhat weaker, than those using simulated data: the corrected median and  
227 maximum-likelihood measures still outperformed all of the five standard measures, while  
228 the corrected mean outperformed three out of the five standard measures (Figure S5a).

229 While the above analysis suggests that the new aggregation measures may be more  
230 accurate than many standard measures over a wide range of conditions, it relied on over  
231 800 estimates to parameterize the individual estimation biases. Such an investment to  
232 characterize estimation biases may be unfeasible for many applications, so we asked how  
233 large of a training dataset is necessary in order to observe improvements in accuracy  
234 over the standard measures. To study this, we obtained a given number of estimates  
235 from across the range of numerosities, generated a maximum-likelihood regression on  
236 that training set, then used that to predict the numerosity of a separate test dataset.  
237 As with the previous analysis, we generated the training and test datasets by drawing  
238 samples using the maximum-likelihood fits shown in Figure 1c-d, but also conducted a  
239 parallel analysis whereby we generated training and test datasets by drawing from our  
240 experimental data (section 3 of the electronic supplementary material for details of both  
241 methodologies).

242 We found rapid improvements in accuracy as the size of the training dataset increased  
243 (Figure 3b). In our simulations, the maximum-likelihood measure begins to outperform  
244 the median and geometric mean when the size of the training dataset is at least 20  
245 samples, the arithmetic mean and trimmed mean after 55 samples, and the average of the  
246 mean and median after 80 samples. The corrected mean required at least 105 samples,  
247 while the corrected median required at least 175 samples, to outperform the five standard  
248 measures. Using samples drawn from our experimental data, our three measures required

249 approximately 200 samples to outperform the five standard measures (Figure S5b). In  
250 short, while our method of correcting biases requires parameterizing bias across the entire  
251 range of numerosities of interest, our simulations show that much fewer training samples  
252 is sufficient for our new aggregation measures to exhibit an accuracy higher than standard  
253 aggregation measures.

254 We next investigated precisely how the size of the test dataset affects accuracy. We  
255 defined an ‘error tolerance’ as the maximum acceptable error of an aggregation measure  
256 and asked what is the probability that a measure achieves a given tolerance for a par-  
257 ticular experiment (the ‘tolerance probability’). As before, we generate test samples by  
258 drawing from the maximum-likelihood fits but also perform an analysis drawing from our  
259 experimental data (see section 4 of the electronic supplementary material for both method-  
260 ologies). For all numerosities, the three new aggregation measures tended to outperform  
261 the five standard measures if the size of the test dataset is relatively large (Figure 4b-c,  
262 Figures S6-S7). However, when the numerosity is large and the size of the test dataset  
263 is relatively small, we observed markedly different patterns. In this regime, the relative  
264 accuracy of aggregation measures can depend on the error tolerance. For example, for  
265 numerosity  $\ln(J) = 10$ , for small error tolerances ( $<0.4$ ), the geometric mean exhibited  
266 the lowest tolerance probability across all of the measures under consideration, but for  
267 large error tolerances ( $>0.75$ ), it is the most likely to fall within tolerance (Figure 4a).  
268 This means that if a researcher wants the collective estimate to be within 40% of the  
269 true value (error tolerance of 0.4), then the geometric mean would be the worst choice  
270 for small test datasets at large numerosities, but if the tolerance was instead set to 75%  
271 of the true value, then the geometric mean would be the best out of all of the measures.  
272 These patterns were also broadly reflected in our analysis using samples drawn from our  
273 experimental data (Figures S8-S10). Therefore, while the corrected measures should have  
274 close to perfect accuracy at the limit of infinite sample size (and perform better than the  
275 standard measures overall), there exist particular regimes in which the standard measures  
276 may outperform the new measures.

277 *3.4. Quantifying the social influence rules*

278 We then conducted an experiment to quantify the social influence rules that individuals  
279 use to update their personal estimate by incorporating information about the estimates  
280 of other people (see Methods for details). Briefly, we first allowed participants to make  
281 an independent estimate. Then we generated artificial ‘social information’ by selecting  
282 a value that was a certain distance from their first estimate (the ‘social distance’), and  
283 informed the participants that this value was the result of averaging across a certain  
284 number of previous estimates (the ‘social group size’). We gave the participants the  
285 opportunity to revise their estimate, and we measured how their change in estimate was  
286 affected by the social distance and social group size. By using artificial information  
287 and masquerading it as real social information, unlike previous studies, we were able to  
288 decouple the effect of social group size, social distance, and the accuracy of the initial  
289 estimate.

290 We found that a fraction of participants (231 out of 602 participants) completely  
291 discounted the social information, meaning that their second estimate was identical to  
292 their first. We constructed a two-stage hurdle model to describe the social influence  
293 rules by first modeling the probability that a participant utilized or discarded social  
294 information, then, for the 371 participants who did utilize social information, we modeled  
295 the magnitude of the effect of social information.

296 A Bayesian approach to fitting a logistic regression model was used to infer whether  
297 social displacement (defined as  $(S - G_1)/G_1$ , where  $S$  is the social estimate and  $G_1$  is  
298 the participant’s initial estimate), social distance (the absolute value of social displace-  
299 ment), or social group size affected the probability that a participant ignored, or used,  
300 social information (see section 5 of the electronic supplementary material for details).  
301 We found that the probability of using social information depends credibly on the social  
302 displacement (coefficient [95% credible interval] = 0.22 [0.03,0.40]), but not on the social  
303 distance (0.061 [-0.12, 0.24]) nor the group size (-0.045 [-0.18, 0.094]) (Figure 5a-c, S11a).  
304 In other words, numerically larger social estimates increased the probability of chang-  
305 ing one’s guess, but numerically smaller social estimates decreased that effect. Posterior

306 predictive checks were used to verify the model captured statistical features of the data  
307 (Figure S12).

308 We next modeled the magnitude of the change in estimate, out of the participants who  
309 did utilize social information. Following [32], we defined a measure of the strength of social  
310 influence,  $\alpha$ , by considering the logarithm of the participant’s revised estimate,  $\ln(G_2)$ ,  
311 as a weighted average of the logarithm of the perceived social information,  $\ln(S)$ , and the  
312 logarithm of the participant’s initial estimate  $\ln(G_1)$ , such that  $\ln(G_2) = \alpha \ln(S) + (1 -$   
313  $\alpha) \ln(G_1)$ . Here,  $\alpha = 0$  indicates that the participant’s two estimates were identical, and  
314 therefore the individual was not influenced by social information at all, while  $\alpha = 1$  means  
315 the participant’s second estimate mirrors the social information. We again used Bayesian  
316 techniques to estimate  $\alpha$  as a normally distributed, logistically transformed linear function  
317 of an intercept, social displacement, social distance, and group size (see section 5 of the  
318 electronic supplementary material for details).

319 Of the subset that changed their estimate, the extent to which they did so depended  
320 credibly on the social displacement (coeff. [95% CI] = 0.65 [0.28, 1.07]), the social distance  
321 (coeff. [95% CI] = -0.41 [-0.82, -0.0052]), and the group size (0.37 [0.17, 0.58]) (Figure  
322 5d-f, S11b). Again, posterior predictive checks revealed the model generated an overall  
323 distribution of social weights consistent with what was found in the data (Figure S13).

### 324 *3.5. The effect of social influence on the wisdom of crowds*

325 If individuals share information with each other before their opinions are aggregated,  
326 then the independent, log-normal distribution of estimates will be altered. Since individu-  
327 als take a form of weighted average of their own estimate and perceived social information,  
328 the distribution of estimates should converge towards intermediate values. However, it is  
329 not clear what effect the observed social influence rules have on the value, or accuracy, of  
330 the aggregation measures [33]. In particular, since the new aggregation measures intro-  
331 duced here were parameterized on independent estimates unaltered by social influence,  
332 their performance may degrade when individuals share information with each other.

333 We simulated several rounds of influence using the rules that we uncovered, using a  
334 social network in which each individual was connected all other individuals, in order to

335 identify measures that may be relatively robust to social influence (see section 6 of the  
336 electronic supplementary material). We used two alternate assumptions about how a set  
337 of estimates is averaged, either by the individual or by an external agent, before being  
338 presented as social information (the ‘individual aggregation measure’), using either the  
339 geometric mean or the arithmetic mean (see section 7 of the electronic supplementary  
340 material). [While the maximum-likelihood measure generally performed the  
341 best in the absence of social influence (Figure 3), this measure was highly  
342 susceptible to the effects of social influence, particularly at large numerosities  
343 (Figure 6). By contrast, the corrected mean was remarkably robust to social  
344 influence, across numerosities and for both individual aggregation measures,  
345 while exhibiting nearly the same accuracy as the maximum-likelihood measure  
346 in the absence of social influence (Figure 3).]

#### 347 4. Discussion

348 While the wisdom of crowds has been documented in many human and non-human  
349 contexts, the limits of its accuracy are still not well understood. Here we demonstrated  
350 how, why, and when collective wisdom may break down by characterizing two major  
351 sources of error, individual (estimation bias) and social (information sharing). We revealed  
352 the limitations of some of the most common averaging measures and introduced three  
353 novel measures that leverage our understanding of these sources of error to improve the  
354 wisdom of crowds.

355 In addition to the conclusions and recommendations drawn for numerosity estima-  
356 tion, the methods described here could be applied to a wide range of estimation tasks.  
357 Estimation biases and social influence are ubiquitous, and estimation tasks may cluster  
358 into broad classes that are prone to similar biases or social rules [34]. For example, the  
359 distribution of estimates for many tasks are likely to be log-normal in nature [35], while  
360 others may tend to be normally distributed. Indeed, there is evidence that counteract-  
361 ing estimation biases can be a successful strategy to improve estimates of probabilities  
362 [36–38], city populations [39], movie box office returns [39], and engineering failure rates

363 [40].

364 Furthermore, the social influence rules that we identified empirically are similar to  
365 general models of social influence, with the exception of the effect of the social displace-  
366 ment that we uncovered. This asymmetric effect suggests that a focal individual was more  
367 strongly affected by social information that was larger in value relative to the focal indi-  
368 vidual’s estimate compared to social information that was smaller than the individual’s  
369 estimate. The observed increase in the coefficient of variation as numerosity increased  
370 (Figure S4b) may suggest that one’s confidence about one’s own estimate decreases as  
371 numerosity increases, which could lead to an asymmetric effect of social distance. Other  
372 estimation contexts in which confidence scales with estimation magnitude could yield a  
373 similar effect. This effect was combined with a weaker negative effect of the social distance,  
374 which is reminiscent of ‘bounded confidence’ opinion dynamics models (*e.g.*, [41–43]),  
375 whereby individuals weight more strongly social information that is similar to their own  
376 opinion. By carefully characterizing both the *individual* estimation biases and *collective*  
377 biases generated by social information sharing, our approach allows us to counteract such  
378 biases, potentially yielding significant improvements when aggregating opinions across  
379 other domains.

380 Other approaches have been used to improve the accuracy of crowds. One strategy  
381 is to search for ‘hidden experts’ and weigh these opinions more strongly [3, 32, 44–47].  
382 While this can be effective in certain contexts, we did not find evidence of hidden experts  
383 in our data. Comparing the group of individuals who ignored social information and those  
384 who utilized social information, the two distribution of estimations were not significantly  
385 different ( $P = 0.938$ , Welch’s t-test on the log-transformed estimates), and the arith-  
386 metic mean, the median, nor our three new aggregation measures were significantly more  
387 accurate across the two groups (Figure S14). In general, searching for hidden experts  
388 requires additional information about the individuals (such as propensity to use social  
389 information, past performance, or confidence level). Our method does not require any  
390 additional information about each individual, only knowledge about statistical tenden-  
391 cies of the population at large (and relatively few samples may be needed to sufficiently



392 parameterize these tendencies).

393 Further refinement of our methods is possible. In cases where the underlying social  
394 network is known [48, 49], or where individuals vary in power or influence [50], simulation  
395 of social influence rules on these networks could lead to a more nuanced understanding of  
396 the mapping between individual to collective estimates. In addition, aggregation measures  
397 can be generalized in a straightforward manner to calculate confidence intervals, in which  
398 an estimate range is generated that includes the true value with some probability. To  
399 improve the accuracy of confidence intervals, information about the sample size and other  
400 features that we showed to be important can be included.

401 In summary, counteracting estimation biases and social influence may be a simple,  
402 general, and computationally efficient strategy to improve the wisdom of crowds.

## 403 **5. Competing interests**

404 We have no competing interests.

## 405 **6. Authors' contributions**

406 ABK, AB, and IDC designed the experiments. ABK, AB, ATH, and MJL performed  
407 the experiments. ABK, AB, JB-C, CCI, and XG analyzed the data. ABK, AB, and IDC  
408 wrote the paper.

## 409 **7. Acknowledgements**

410 We thank Stefan Krause, Jens Krause, Andrew King, and Michael J. Mauboussin  
411 for contributing datasets to this study, and Mirta Galesic for providing feedback on the  
412 manuscript.

## 413 **8. Data Accessibility**

414 Datasets are available in the electronic supplementary material.

415 **9. Ethics**

416 The experimental procedures were approved by the Princeton University and Santa  
417 Fe Institute ethics committees.

418 **10. Funding**

419 ABK was supported by a James S. McDonnell Foundation Postdoctoral Fellowship  
420 Award in Studying Complex Systems. AB was supported by an SFI Omidyar Post-  
421 doctoral Fellowship. CCI was supported by a NERC Independent Research Fellowship  
422 NE/K009370/1. IDC acknowledges support from NSF (PHY-0848755, IOS-1355061,  
423 EAGER-IOS-1251585), ONR (N00014-09-1-1074, N00014-14-1-0635), ARO (W911NG-  
424 11-1-0385, W911NF-14-1-0431), and the Human Frontier Science Program (RGP0065/2012).

- 425 [1] J. Surowiecki, *The Wisdom of the Crowds: Why the Many are Smarter than the*  
426 *Few*, Little Brown, 2004.
- 427 [2] F. Galton, *Vox populi*, *Nature* 75 (1907) 450–451.
- 428 [3] D. Prelec, H. Seung, J. McCoy, A solution to the single-question crowd wisdom  
429 problem, *Nature* 541 (2017) 532–535.
- 430 [4] B. Bahrami, K. Olsen, P. Latham, A. Roepstorff, G. Rees, C. Frith, Optimally  
431 interacting minds, *Science* 329 (2010) 1081–1085.
- 432 [5] S. Krause, R. James, J. Faria, G. Ruxton, J. Krause, Swarm intelligence in humans:  
433 diversity can trump ability, *Animal Behaviour* 81 (2011) 941–948.
- 434 [6] A. King, L. Cheng, S. Starke, J. Myatt, Is the true ‘wisdom of the crowd’ to copy  
435 successful individuals?, *Biology Letters* 8 (2011) 197–200.
- 436 [7] D. Sumpter, J. Krause, R. James, I. Couzin, A. Ward, Consensus decision making  
437 by fish, *Current Biology* 18 (2008) 1773–1777.
- 438 [8] A. Ward, J. Herbert-Read, D. Sumpter, J. Krause, Fast and accurate decisions  
439 through collective vigilance in fish shoals, *Proc Natl Acad Sci USA* 108 (2011) 2312–  
440 2315.
- 441 [9] A. Ward, D. Sumpter, I. Couzin, P. Hart, J. Krause, Quorum decision-making  
442 facilitates information transfer in fish shoals, *Proc Natl Acad Sci USA* 105 (2008)  
443 6948–6953.
- 444 [10] T. Sasaki, B. Granovskiy, R. Mann, D. Sumpter, S. Pratt, Ant colonies outperform  
445 individuals when a sensory discrimination task is difficult but not when it is easy,  
446 *Proc Natl Acad Sci USA* 110 (2013) 13769–13773.
- 447 [11] T. Sasaki, S. Pratt, Emergence of group rationality from irrational individuals, *Behav*  
448 *Ecol* 22 (2011) 276–281.

- 449 [12] S. Tamm, Bird orientation: single homing pigeons compared to small flocks, *Behav*  
450 *Ecol Sociobiol* 7 (1980) 319–322.
- 451 [13] A. Simons, Many wrongs: the advantage of group navigation, *Trends Ecol Evol* 19  
452 (2004) 453–455.
- 453 [14] N. Condorcet, *Essai sur l’application de l’analyse à la probabilité des décisions ren-*  
454 *dues à la pluralité de voix*, Paris, 1785.
- 455 [15] D. Kahneman, *Thinking, Fast and Slow*, Straus and Giroux, 2011.
- 456 [16] R. Nickerson, Confirmation bias: a ubiquitous phenomenon in many guises, *Review*  
457 *of General Psychology* 2 (1998) 175–220.
- 458 [17] M. Haselton, D. Nettle, The paranoid optimist: an integrative evolutionary model  
459 of cognitive biases, *Personality and Social Psychology Review* 10 (2006) 47–66.
- 460 [18] S. Stevens, On the psychophysical law, *The Psychological Review* 64 (1957) 153–181.
- 461 [19] J. Whalen, C. Gallistel, R. Gelman, Nonverbal counting in humans: the psy-  
462 chophysics of number representation, *Psychological Science* 10 (1999) 130–137.
- 463 [20] V. Izard, S. Dehaene, Calibrating the mental number line, *Cognition* 106 (2008)  
464 1221–1247.
- 465 [21] L. E. Krueger, Single judgments of numerosity, *Attention, Perception, & Psy-*  
466 *chophysics* 31 (1982) 175–182.
- 467 [22] L. E. Krueger, Perceived numerosity: A comparison of magnitude production, mag-  
468 nitude estimation, and discrimination judgments, *Attention, Perception, & Psy-*  
469 *chophysics* 35 (1984) 536–542.
- 470 [23] L. Krueger, Reconciling fechner and stevens: toward a unified psychophysical law,  
471 *Behavioral and Brain Sciences* 12 (1989) 251–320.
- 472 [24] F. Galton, One vote, one value, *Nature* 75 (1907) 414.

- 473 [25] J. Lorenz, H. Rauhut, F. Schweitzer, D. Helbing, How social influence can undermine  
474 the wisdom of crowd effect, *Proc Natl Acad Sci USA* 108 (2011) 9020–9025.
- 475 [26] J. Armstrong, Combining forecasts. Armstrong, JS, ed. *Principles of Forecasting: A*  
476 *Handbook for Researchers and Practitioners*, Kluwer, New York, 2001.
- 477 [27] M. Lobo, D. Yao, Human judgement is heavy tailed: Empirical evidence and im-  
478 plications for the aggregation of estimates and forecasts, Fontainebleau: INSEAD,  
479 2010.
- 480 [28] A. Kao, N. Miller, C. Torney, A. Hartnett, I. Couzin, Collective learning and optimal  
481 consensus decisions in social animal groups, *PLoS Computational Biology* 10 (2014)  
482 e1003762.
- 483 [29] C. Wagner, C. Schneider, S. Zhao, H. Chen, The wisdom of reluctant crowds, Pro-  
484 ceedings of the 43rd Hawaii International Conference on System Sciences, 2010.
- 485 [30] M. Mauboussin, Explaining the wisdom of crowds, Legg Mason Capital Management  
486 White Paper, 2007.
- 487 [31] S. Kemp, Estimating the sizes of sports crowds, *Perceptual and motor skills* 59  
488 (1984) 723–729.
- 489 [32] G. Madirolas, G. de Polavieja, Improving collective estimations using resistance to  
490 social influence, *PLoS Comput Biol* 11 (2015) e1004594.
- 491 [33] B. Golub, M. Jackson, Naïve learning in social networks and the wisdom of crowds,  
492 *American Economic Journal: Microeconomics* 2 (2010) 112–149.
- 493 [34] M. Steyvers, B. Miller, Cognition and collective intelligence, *Handbook of Collective*  
494 *Intelligence* (2015) 119.
- 495 [35] S. Dehaene, V. Izard, E. Spelke, P. Pica, Log or linear? distinct intuitions of the  
496 number scale in western and amazonian indigene cultures, *Science* 320 (2008) 1217–  
497 1220.

- 498 [36] B. M. Turner, M. Steyvers, E. C. Merkle, D. V. Budescu, T. S. Wallsten, Forecast  
499 aggregation via recalibration, *Machine learning* 95 (2014) 261–289.
- 500 [37] M. D. Lee, I. Danileiko, Using cognitive models to combine probability estimates,  
501 *Judgment and Decision Making* 9 (2014) 259.
- 502 [38] V. A. Satopää, J. Baron, D. P. Foster, B. A. Mellers, P. E. Tetlock, L. H. Ungar,  
503 Combining multiple probability predictions using a simple logit model, *International*  
504 *Journal of Forecasting* 30 (2014) 344–356.
- 505 [39] A. Whalen, S. Yeung, Using ground truths to improve wisdom of the crowd estimates,  
506 *Proceedings of the Annual Cognitive Science Society Meeting* (2015).
- 507 [40] E. Merkle, M. Steyvers, A psychological model for aggregating judgments of magni-  
508 tude, *Social computing, behavioral-cultural modeling and prediction* (2011) 236–243.
- 509 [41] R. Hegselmann, U. Krause, Opinion dynamics and bounded confidence models,  
510 analysis, and simulation, *J Artif Soc Soc Simulat* 5 (2002).
- 511 [42] G. Deffuant, D. Neau, F. Amblard, G. Weisbuch, Mixing beliefs among interacting  
512 agents, *Advances in Complex Systems* 3 (2000) 87–98.
- 513 [43] G. Deffuant, F. Amblard, G. Weisbuch, T. Faure, How can extremism prevail? a  
514 study based on the relative agreement interaction model, *Journal of artificial societies*  
515 *and social simulation* 5 (2002).
- 516 [44] A. Koriat, When are two heads better than one and why?, *Science* 336 (2012)  
517 360–362.
- 518 [45] S. Hill, N. Ready-Campbell, Expert stock picker: the wisdom of (experts in) crowds,  
519 *International Journal of Electronic Commerce* 15 (2011) 73–102.
- 520 [46] J. Whitehill, T. Wu, J. Bergsma, J. Movellan, P. Ruvolo, Whose vote should count  
521 more: optimal integration of labels from labelers of unknown expertise, *Advances in*  
522 *Neural Information Processing Systems* 22 (2009) 2035–2043.

- 523 [47] D. V. Budescu, E. Chen, Identifying expertise to extract the wisdom of crowds,  
524 Management Science 61 (2014) 267–280.
- 525 [48] M. Jönsson, U. Hahn, E. Olsson, The kind of group you want to belong to: Effects  
526 of group structure on group accuracy, Cognition 142 (2015) 191–204.
- 527 [49] M. Moussaïd, S. Herzog, J. Kämmer, R. Hertwig, Reach and speed of judgment  
528 propagation in the laboratory, Proc Natl Acad Sci USA 114 (2017) 4117–4122.
- 529 [50] J. Becker, D. Brackbill, D. Centola, Network dynamics of social influence in the  
530 wisdom of crowds, Proc Natl Acad Sci USA 114 (2017) E5070–E5076.

531 **Figure 1. The effect of numerosity on the distribution of estimates.** (a) An  
532 example jar containing 659 objects ( $\ln(J) = 6.5$ ). (b) The histogram of estimates (grey  
533 bars) resulting from the jar shown in (a) closely approximates a log-normal distribution  
534 (solid black line); dotted vertical line indicates the true number of objects. A log-normal  
535 distribution is described by two parameters,  $\mu$  and  $\sigma$ , which are the mean and standard  
536 deviation, respectively, of the normal distribution that results when the logarithm of the  
537 estimates is taken (inset). (c-d) The two parameters  $\mu$  and  $\sigma$  increase linearly with the  
538 logarithm of the true number of objects,  $\ln(J)$ . Solid lines: maximum-likelihood estimate,  
539 shaded area: 95% confidence interval. The maximum-likelihood estimate was calculated  
540 using only the five original datasets collected for this study (black circles); the eight other  
541 datasets collected from the literature are shown only for comparison (grey circles indicate  
542 other datasets for which the full dataset was available, white circles indicate datasets for  
543 which only summary statistics were available, see section 1 of the electronic supplementary  
544 material).

545 **Figure 2. The accuracy of the arithmetic mean and the median.** (a) The  
546 expected value of the arithmetic mean (blue) and median (red) relative to the true number  
547 of objects (black dotted line), as a function of  $\ln(J)$ . The relative value is defined as  
548  $(X - J)/J$ , where  $X$  is the value of the aggregation measure. (b) The relative error of the  
549 expected value of the two aggregation measures, defined as  $|X - J|/J$ . For both panels,  
550 solid lines indicate maximum-likelihood values, shaded areas indicate 95% confidence  
551 intervals, and solid circles show the empirical values from the five datasets.

552 **Figure 3. The overall relative performance of the aggregation measures.**  
553 (a) The percentage of simulations in which the measure indicated in the row was more  
554 accurate than the measure indicated in the column. The three new measures are listed  
555 in the rows and are compared to all eight measures in the columns. Colors correlate with  
556 percentages (blue:  $>50\%$ , red:  $<50\%$ ). (b) The median error of the three new aggregation  
557 measures (corrected median, dashed red line; corrected mean, dashed blue line; maximum-



558 likelihood measure, dashed green line) as a function of the size of the training dataset. The  
559 three new aggregation measures are compared against the arithmetic mean (solid blue),  
560 median (solid red), the geometric mean (orange), the average of the mean and the median  
561 (yellow), and the trimmed mean (magenta). The 95% confidence interval are displayed  
562 for the latter measures, which are not a function of the size of the training dataset.

563 **Figure 4. The effect of the test dataset size and error tolerance level on the**  
564 **relative accuracy of the aggregation measures.** The probability that an aggregation  
565 measure exhibits a relative error (defined as  $|X - J|/J$ , where  $X$  is the value of an  
566 aggregation measure) less than a given error tolerance, for test dataset size (a) 4, (b)  
567 64, and (c) 512, and numerosity  $J = 22026$  ( $\ln(J) = 10$ ). In panel (a), the lines for the  
568 arithmetic mean and the trimmed mean are nearly identical; in panel (c), the lines for  
569 the corrected mean and corrected median are nearly identical.

570 **Figure 5. The social influence rules.** [The probability that an individual is  
571 affected by social information (a) increases with social displacement (the rel-  
572 ative distance between the value of a participant's estimate and the value of  
573 the social information) but does not depend on (b) the social distance (the  
574 absolute distance between a participant's estimate and the social information)  
575 or (c) the perceived social group size. The social influence weight  $\alpha$  (d) in-  
576 creases with social displacement, (e) decreases with social distance, and (f)  
577 increases with social group size. Solid lines: predicted mean value; shaded  
578 area: 95% credible interval; circles: the mean of binned data for (a-c) and raw  
579 data for (d-f); red lines and areas indicate a credible effect (see Figure S13  
580 for the posterior distributions of each coefficient). We note that some of the  
581 empirical data extend outside of the bounds of the plots in (d-f); we selected  
582 the bounds to more clearly show the patterns of the fitted parameters.]

583 **Figure 6. The robustness of aggregation measures under social influence.** The  
584 relative error of the eight aggregation measures without social influence (gray circles) and  
585 after ten rounds of social influence (black circles) when (a-c) individuals internally take  
586 the geometric mean of the social information that they observe, or when (d-f) individuals  
587 internally take the arithmetic mean of the social information, for numerosity  $\ln(J) = 4$   
588 (a,d),  $\ln(J) = 7$  (b,e), and  $\ln(J) = 10$  (c,f). Circles show the mean relative error across  
589 1000 replicates, error bars show twice the standard error. The error bars are often smaller  
590 than the size of the corresponding circles, and where some gray circles are not visible,  
591 they are nearly identical to the corresponding black circles.