

# Surprisingly Popular Voting Recovers Rankings, Surprisingly!

Hadi Hosseini<sup>1</sup>, Debmalya Mandal<sup>2</sup>, Nisarg Shah<sup>3</sup> and Kevin Shi<sup>3</sup>

<sup>1</sup>Pennsylvania State University

<sup>2</sup>Columbia University

<sup>3</sup>University of Toronto

hadi@psu.edu, dm3557@columbia.edu, nisarg@cs.toronto.edu, kevins.shi@mail.utoronto.ca

## Abstract

The wisdom of the crowd has long become the de facto approach for eliciting information from individuals or experts in order to predict the ground truth. However, classical democratic approaches for aggregating individual *votes* only work when the opinion of the majority of the crowd is relatively accurate. A clever recent approach, *surprisingly popular voting*, elicits additional information from the individuals, namely their *prediction* of other individuals' votes, and provably recovers the ground truth even when experts are in minority. This approach works well when the goal is to pick the correct option from a small list, but when the goal is to recover a true ranking of the alternatives, a direct application of the approach requires eliciting too much information. We explore practical techniques for extending the surprisingly popular algorithm to ranked voting by partial votes and predictions and designing robust aggregation rules. We experimentally demonstrate that even a little prediction information helps surprisingly popular voting outperform classical approaches.

## 1 Introduction

The wisdom of the crowd has been the default choice for uncovering the ground truth. Suppose we wish to determine the true answer to the question: "Is Philadelphia the capital of Pennsylvania?" Condorcet's Jury Theorem suggests that if we elicit votes from a large crowd, the majority answer will be correct with high probability even if, on average, the crowd is only slightly more accurate than a random selection. However, in some domains the crowd can be highly inaccurate and experts may be in minority. For example, when the very question listed above is posed to real crowds, the majority answer is often (the incorrect) 'yes' [De Boer and Bernstein, 2017].

To circumvent this difficulty and uncover the ground truth even when the majority is wrong, Prelec *et al.* [2017] introduce the *surprisingly popular* (SP) algorithm. This algorithm asks each individual not only what she thinks the answer is (the *vote*), but also what fraction of the other participants she thinks will say yes/no (the *prediction*). Then, instead of simply selecting the majority (i.e. popular) answer, the algorithm

selects the answer that is *surprisingly popular*, i.e., whose actual frequency in the votes is greater than its average predicted frequency. They show that as the crowd gets larger in the limit, this approach will provably recover the correct answer with probability 1, even if the crowd is less accurate than a random selection on average.

The intuition behind their algorithm, borrowed from their work, is as follows. Suppose there are two hypothetical worlds, one where Philadelphia is the capital and one where it is not. In the former world, a greater fraction (say 90%) would say 'yes' than the fraction (say 60%) that would say 'yes' in the latter. However, the 60% of the people who believe the correct world is the former would predict the frequency of 'yes' to be 90%, whereas the remaining 40% would predict it to be 60%. This would make the average predicted frequency of 'yes' to be somewhere between 60% and 90%, higher than its actual frequency of 60%. In other words, the majority but incorrect answer 'yes' would be surprisingly *unpopular* while 'no' would be surprisingly popular and correct.

Several works have demonstrated the effectiveness of this approach in a wide range of domains [Prelec *et al.*, 2017; Lee *et al.*, 2018; Wang *et al.*, 2019; Palley and Soll, 2019; Rutchick *et al.*, 2020; Mandal *et al.*, 2020a]. Prediction questions have also been used to boost the accuracy of surveys on social networks [Galesic *et al.*, 2018]. Prelec *et al.* [2017] show how to apply their approach to questions with non-binary votes and non-binary ground truth. When the true answer lurks among  $r$  options, their approach requires each individual to predict the exact frequency of each of  $r$  options among other individuals' votes. We are interested in ranked voting, i.e., when the ground truth is a ranking of  $m$  alternatives. Note that in this case, the approach of Prelec *et al.* [2017], which we refer to as *surprisingly popular* (SP) *voting*, would require eliciting predictions in the form of a distribution over  $r = m!$  options, which is clearly infeasible for even moderate values of  $m$ . Thus, the main research questions we address are:

*How do we extend surprisingly popular voting to effectively recover a ground truth ranking of alternatives? If we elicit partial vote and prediction, how do we aggregate them and what information-accuracy tradeoff does this offer?*

## 1.1 Our Contributions

We focus on eliciting only *ordinal* vote and prediction information. For the *vote*, we ask individuals to provide their opinion of either just the top alternative of the ground truth ranking (*Top*) or the full ground truth ranking (*Rank*). For the *prediction*, informally, we ask individuals to predict either just a single alternative (*Top*) or a ranking of alternatives (*Rank*) based on the other individuals' votes. The exact prediction elicited under various conditions is described in Section 3. In addition to these four elicitation formats, we use as benchmark two classical elicitation formats in which Top and Rank votes are elicited but no prediction is elicited. Because the SP algorithm of Prelec *et al.* [2017] does not work on partial votes and predictions, we first design a novel aggregation method for such partial information.

Next, we conduct an empirical study with 720 participants from Amazon's Mechanical Turk platform. We ask the participants questions on geography, movies, and artwork which admit a ground truth ranking of four alternatives and elicit their responses in the aforementioned six elicitation formats. We compare the different elicitation formats using four metrics: difficulty (measured through response time as well as perceived difficulty), expressiveness, error in recovering the ground truth top alternative, and error in recovering the ground truth ranking.

Our results show that even when the vote and prediction information are individually no better than random guesses, by combining the two pieces of information SP voting performs significantly better. Further, it outperforms a whole slew of conventional voting rules which ignore prediction information and only aggregate the votes. We also observe that when it is necessary to choose between eliciting more complex vote information and eliciting more complex prediction information, the latter may be the right choice.

## 1.2 Related Work

Our work builds on the SP voting approach of Prelec *et al.* [2017]. This approach in turn builds on its precursor, the Bayesian truth serum (BTS) [Prelec, 2004], which also uses participants' predictions, but for a different objective: to decide payoffs to the participants which incentivize them to honestly report their votes and predictions.

Prediction markets [Arrow *et al.*, 2008; Chen and Pennock, 2010], quadratic voting [Lalley and Weyl, 2018], and peer prediction [Miller *et al.*, 2005] are alternative approaches to recovering the ground truth, which, like SP voting, allow a minority of experts to override the majority opinion. Instead of eliciting participants' predictions of other participants' votes, prediction markets and quadratic voting ask participants to place a bet on their vote while peer prediction methods require them to participate in multiple tasks.

These recent approaches stand in contrast to a large body of work on epistemic social choice [Pivato, 2019] and noisy voting [Caragiannis *et al.*, 2016], which build on the seminal work of de Condorcet [1785], Galton [1907], and Young [1988]. Some of this literature focuses on statistical models of errors in participants' votes such as the Mallows model, the Bradley-Terry model, the Thurstone-Mosteller

model, and the Plackett-Luce model. However, all these models assume that a participant is ever-so-slightly more likely to report the correct option than an incorrect option. Hence, approaches based on these models can fail to recover the ground truth when the majority of the crowd is misinformed.

Finally, our work is reminiscent of a recent flurry of work on the elicitation-distortion tradeoff in computational social choice [Mandal *et al.*, 2019; Abramowitz *et al.*, 2019; Mandal *et al.*, 2020b; Kempe, 2020; Amanatidis *et al.*, 2020]. In this line of work, there is no ground truth; instead, participants have subjective preferences and the goal is to identify the decision that maximizes the social welfare. Rather than directly eliciting participants' utility functions, various elicitation formats are used to elicit partial preferences to analyze the tradeoff between the amount of information elicited and the approximation to social welfare (called distortion). Our work replaces the distortion with its counterpart, that is, the accuracy of recovering an underlying ground truth.

## 2 Model

Let  $A$  be a set of  $m$  alternatives and  $\mathcal{L}(A)$  be the set of rankings over  $A$ . For a ranking  $\sigma \in \mathcal{L}(A)$  and  $x \in \{1, \dots, m\}$ , let  $\sigma(x)$  be the alternative in the  $x^{\text{th}}$  highest position in  $\sigma$ .

SP voting uses a Bayesian model; in the following, we present a special case of the model for ranked voting. There exists a ground truth ranking  $\pi^* \in \mathcal{L}(A)$  drawn from a *prior*  $\mathcal{P}$ . There are  $n$  voters; each voter  $i$  observes a noisy ranking  $\sigma_i \in \mathcal{L}(A)$  drawn from a *signal distribution*  $\Pr_s(\cdot|\pi^*)$ . The voters know both the prior  $\mathcal{P}$  and the signal distribution  $\Pr_s(\cdot|\pi^*)$ ; however, the principal is unaware of both. Following Prelec *et al.* [2017], we assume that  $\mathcal{P}(\pi), \Pr_s(\sigma|\pi) > 0$  for all rankings  $\sigma, \pi \in \mathcal{L}(A)$  to avoid degeneracy.

Conventional voting would ask each voter  $i$  to simply report her observed noisy ranking  $\sigma_i$  and use a voting rule such as the Kemeny rule or Borda count to aggregate the reported rankings. SP voting additionally asks each voter  $i$  to make inferences about the reports of other voters. Given her observed noisy ranking  $\sigma_i$  and the prior  $\mathcal{P}$ , voter  $i$  can compute a posterior distribution over the ground truth, given by

$$\Pr_g(\pi^*|\sigma_i) = \frac{\Pr_s(\sigma_i|\pi^*) \cdot \mathcal{P}(\pi^*)}{\sum_{\pi' \in \mathcal{L}(A)} \Pr_s(\sigma_i|\pi') \cdot \mathcal{P}(\pi')}.$$

In turn, the voter can also infer a distribution over the noisy ranking  $\sigma_j$  observed by another voter  $j$ :

$$\Pr_o(\sigma_j|\sigma_i) = \sum_{\pi^* \in \mathcal{L}(A)} \Pr_s(\sigma_j|\pi^*) \cdot \Pr_g(\pi^*|\sigma_i).$$

SP voting asks each voter  $i$  to report not only her observed noisy ranking  $\sigma_i$  (the *vote*), but also her inferred distribution  $\Pr_o(\cdot|\sigma_i)$  over other voters' noisy rankings (the *prediction*). Given these reports, for a ranking  $\pi \in \mathcal{L}(A)$ , let  $f(\pi) = \sum_{i=1}^n \mathbb{1}[\sigma_i = \pi]$  denote the number of voters who vote  $\pi$  and  $g(\cdot|\pi)$  denote the average of reported predictions  $\Pr_o(\cdot|\sigma_i)$  across all voters  $i$  with  $\sigma_i = \pi$ . Then, the SP algorithm of Prelec *et al.* [2017] computes the prediction-normalized vote count for each possible ground truth  $\pi$  as

$$\bar{V}(\pi) = f(\pi) \cdot \sum_{\pi' \in \mathcal{L}(A)} \frac{g(\pi'|\pi)}{g(\pi|\pi')}. \quad (1)$$

The following result due to Prelec *et al.* [2017], rephrased in our context, guarantees that the ground truth ranking will have the highest prediction-normalized vote count under the assumption that the highest posterior probability for ground truth ranking  $\pi$  will be assigned by a voter who observes noisy ranking  $\pi$ .

**Theorem 1** ([Prelec *et al.*, 2017]). *Suppose the prior  $\mathcal{P}$  and the signal distribution  $\text{Pr}_s$  are such that  $\text{Pr}_g(\pi|\pi) > \text{Pr}_g(\pi|\pi')$  for all distinct rankings  $\pi, \pi' \in \mathcal{L}(A)$ . Then, we have that  $\text{Pr}[\pi^* \in \text{argmax}_{\pi \in \mathcal{L}(A)} \bar{V}(\pi)] \rightarrow 1$  as  $n \rightarrow \infty$ .*

### 3 Elicitation Formats & Aggregation Rules

Note that the prediction requested from voter  $i$ ,  $\text{Pr}_o(\cdot|\sigma_i)$ , is a distribution over  $m!$  rankings. Eliciting this would undoubtedly place significant cognitive burden on the voter. Thus, our goal is to elicit partial vote and prediction information from the voters. Since eliciting numerical information is known to be difficult [Camerer, 2011], we focus on eliciting ordinal information for prediction. We develop aggregation rules for recovering the ground truth from ordinal information and empirically evaluate the effectiveness of SP voting.

#### 3.1 Elicitation Formats

We focus on two types of vote reports, and for each of them, two types of prediction reports. Below we provide formal explanations of these formats in the context of our model. In the next section, we provide example phrasings that were used to pose the various questions to the participants in our empirical study. Let  $r_i$  and  $q_i$  respectively denote the vote and prediction reports submitted by voter  $i$ .

- *Top vote:* Voter  $i$  reports the top alternative in her observed noisy ranking, i.e.,  $r_i = \sigma_i(1)$ .
  - *Top prediction:* Voter  $i$  estimates the most frequent alternative among the other votes, i.e.  $q_i = \text{argmax}_{a \in A} \sum_{\sigma \in \mathcal{L}(A): \sigma(1)=a} \text{Pr}_o(\sigma|\sigma_i)$ .
  - *Rank prediction:* Voter  $i$  estimates the ranking of the alternatives by their frequency among the other votes, i.e.,  $q_i \in \mathcal{L}(A)$  such that  $\sum_{\sigma \in \mathcal{L}(A): \sigma(1)=q_i(x)} \text{Pr}_o(\sigma|\sigma_i) \geq \sum_{\sigma \in \mathcal{L}(A): \sigma(1)=q_i(y)} \text{Pr}_o(\sigma|\sigma_i)$  for all  $x > y$ .
- *Rank vote:* Voter  $i$  reports her entire observed noisy ranking, i.e.,  $r_i = \sigma_i$ .
  - *Top prediction:* Voter  $i$  estimates the alternative that appears most frequently in the top position of the other votes. Formally, this is equivalent to the top prediction in case of a top vote:  $q_i = \text{argmax}_{a \in A} \sum_{\sigma \in \mathcal{L}(A): \sigma(1)=a} \text{Pr}_o(\sigma|\sigma_i)$ .
  - *Rank prediction:* Voter  $i$  estimates the most frequent ranking among the other votes, i.e.,  $q_i \in \text{argmax}_{\sigma \in \mathcal{L}(A)} \text{Pr}_o(\sigma|\sigma_i)$ . Note that this is different from the rank prediction in case of a top vote.

This gives rise to four elicitation formats, which we refer to as Top-Top, Top-Rank, Rank-Top, and Rank-Rank with the first component denoting the vote format and the second

denoting the prediction format. As a benchmark, we use Top-None and Rank-None, where top and rank votes are elicited, respectively, but no prediction information is elicited.

#### 3.2 Aggregation Rules

There are two difficulties in applying the SP algorithm of Prelec *et al.* [2017] — maximizing  $\bar{V}(\pi)$  given in Equation (1) — in our setting.

First, the effectiveness of the approach depends on how accurately functions  $f$  and  $g$  from Equation (1) match their expected values, which in turn depends on how large the number of voters is compared to the number of options among which the ground truth lurks. In our case, since the ground truth is one of  $m!$  rankings, the approach would be ineffective unless each question is answered by a number of voters much larger than  $m!$ . Instead, we determine the ground truth comparison of each of  $\binom{m}{2}$  pairs of alternatives independently by applying the algorithm from Equation (1) on the relevant pairwise comparison data extracted from the reports of the voters.

Second, even for comparing a pair of alternatives, Equation (1) requires cardinal prediction information whereas our input is ordinal. We propose a simple parametric model in which, for each elicitation format, we use two parameters,  $\alpha \in (0.5, 1)$  and  $\beta \in (0, 0.5)$ , to convert ordinal pairwise predictions into cardinal pairwise predictions to be utilized by the SP algorithm. In Section 4, we describe how we train these parameter values. The formal algorithm and its detailed description are provided in the full version.<sup>1</sup>

Note that applying our algorithm for comparing each pair of alternatives independently results in a tournament, which we use for two prediction tasks: predicting the top alternative in the ground truth ranking and predicting the entire ground truth ranking. For the former task, we select the alternative that defeats the maximum number of other alternatives in the resulting tournament, breaking ties uniformly at random, and consider the frequency of predicting the correct top alternative. For the latter task, we compute the Kendall Tau distance of the tournament from the ground truth ranking.

Finally, note that there are no prediction reports for Top-None and Rank-None and we consider a natural extension of SP voting. In particular, for Top-None, SP voting returns an acyclic tournament comparing alternatives by their plurality scores, and for Rank-None, it returns the (potentially cyclic) majority preference tournament. We then select an alternative/ranking as described earlier.

### 4 Experiment Design

To test the effectiveness of SP voting for recovering ranked ground truth with only ordinal elicitation, we conducted an empirical study by recruiting 720 participants (turkers) from Amazon Mechanical Turk (MTurk), a popular crowdsourcing marketplace. An average turker spent about 15 minutes to complete the survey. The survey was designed as follows.

**Datasets.** To generate questions with an underlying ground truth comparison of alternatives, we used three datasets from three distinct *domains*:

<sup>1</sup><https://arxiv.org/abs/2105.09386>

1. The *geography* dataset<sup>2</sup> contains 230 countries with their 2019 population estimates according to the United Nations.
2. The *movies* dataset<sup>3</sup> contains 15,743 movies with their lifetime box-office gross earnings.
3. The *paintings* dataset<sup>4</sup> contains 80 paintings with their latest auction prices.

**Questions.** In each domain, the numerical values associated with the alternatives allow a ground truth comparison among the alternatives. For each domain, we considered the top 50 alternatives with the highest values. From these, we generated 20 questions, each comparing four alternatives selected such that two consecutive alternatives in the ground truth ranking were exactly 6 ranks apart in the global ranking of all 50 alternatives. Collectively, we had 60 questions across all three domains. For each of the 60 questions and each of the 6 elicitation formats described in Section 3, we elicited 20 responses, generating a total of 7,200 responses.

**Turker Assignment.** Figure 1 shows the workflow faced by a turker. Each of the 720 turkers responded to 10 questions split evenly among two randomly assigned elicitation formats. The turkers were divided roughly equally between the 30 ordered pairs of elicitation formats called *treatments*. Further, as mentioned above, each question under each elicitation format was assigned to the same number of turkers.

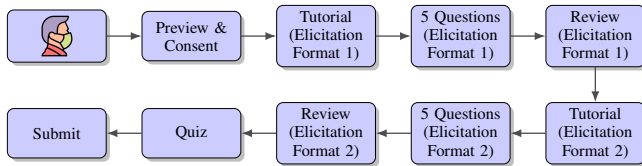


Figure 1: The workflow of a turker.

**Tutorials.** As shown in Figure 1, each set of five questions in a fixed elicitation format was preceded by a tutorial. The tutorial was designed specifically for the elicitation format and tested turkers’ understanding of the vote and prediction formats. It contained a sample question along with pre-specified beliefs over the correct answer as well as over the other responses. Turkers had to successfully pass the tutorial by converting the given beliefs into the requested vote and prediction format in order to proceed to the questions.

**Reviews.** Each set of five questions was also succeeded by a review, which asked the turkers to rate the *difficulty* (from Very Easy to Very Difficult) and *expressiveness* (Very Little to Very Significant) of the elicitation format of the preceding questions. While we controlled the difficulty level of various questions from a given domain, as we show in Section 5 the three domains themselves differed significantly in their difficulty. In anticipation of this and to ensure that the turkers’ implicit comparison between their two assigned elicitation formats is not influenced by the domains, the study was

designed such that the sequence of domains encountered by a turker in the first five questions precisely matched that in the next five questions. See the full version for details such as the consent form, the tutorial for each domain, the review, and other details.

**Response Qualifications.** To ensure high-quality responses, in addition to providing training in the form of tutorials, we restricted participation in our study to turkers who had (a) at least 90% approval rate on previous tasks, (b) at least 100 completed tasks, and (c) the region set to US East (us-east-1) on MTurk. Additionally, at the end of the survey, the turkers were required to answer a quiz, which repeated the four alternatives from the last question they answered and asked them to identify the alternative they chose or ranked first in their vote. The turkers were incentivized to answer the quiz correctly (see below). In our case, over 82% of turkers passed the quiz.

**Payments.** The payment was divided into two parts. A *base* payment of 50¢ was provided conditioned on completing the entire survey including all tutorials, questions, and reviews. A *bonus* payment of 50¢ was provided conditioned on correctly answering the quiz question.

**Elicitation Formats.** In Section 3, we discussed six elicitation formats and described what vote and prediction a given voter  $i$  should submit as a function of her observed noisy ranking  $\sigma_i$ , the prior  $\mathcal{P}$ , and the signal distribution  $\text{Pr}_s$ . In our empirical study, we design natural and intuitive phrasing to elicit the corresponding responses from the turkers. The full version of the paper contains sample phrasings for all six elicitation formats and screenshots from our user interface. Here we give one example for the Top-Rank elicitation format. Consider a question which asks to compare four countries (United Kingdom, Vietnam, Russia, and Kenya) by their population. Under the Top-Rank elicitation format, the vote and prediction questions would be as follows:

- **Part A (vote):** Which country do you think is the most populated among the following?
- **Part B (prediction):** Imagine that other participants will also answer Part A. How do you think the following countries will be ordered from the most common response (top) to the least common (bottom)?

**Training.** Recall that in our aggregation method, for each elicitation format, we use two parameters,  $\alpha \in (0.5, 1)$  and  $\beta \in (0, 0.5)$ , to convert ordinal predictions into cardinal predictions that can be then used in the SP algorithm. To learn effective values of these parameters, we split the dataset into a training and a test set. For each elicitation format, we selected 5 questions from each of three domains, reserving the remaining 15 questions from each domain for the test set. Using these 15 questions, we performed a grid search over  $\alpha$  ranging from 0.55 to 0.95 in increments of 0.025 and  $\beta$  ranging from 0.05 to 0.45 in increments of 0.025 and selected the values with the lowest mean squared error.

## 5 Results

In this section, we present our results averaged across all three domains. In the full version, we present more detailed results

<sup>2</sup>Retrieved from worldpopulationreview.com

<sup>3</sup>Retrieved from boxofficemojo.com/chart/top\_lifetime\_gross

<sup>4</sup>Generously provided by the authors of Prelec *et al.* [2017].

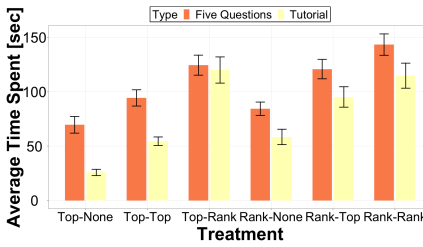


Figure 2: Average time spent.

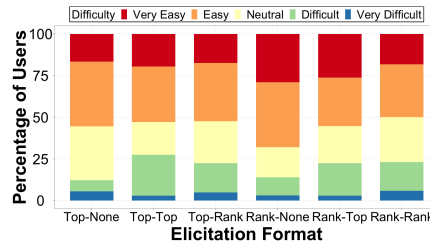


Figure 3: Perceived difficulty.

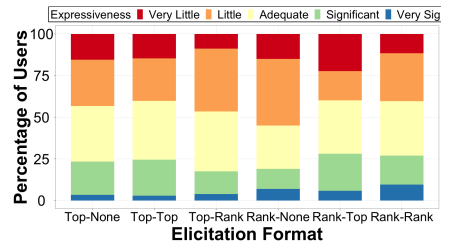


Figure 4: Perceived expressiveness.

averaged across each domain separately. All confidence intervals shown are 95% intervals. We compare the elicitation formats using four key metrics: difficulty (i.e. cognitive burden), expressiveness, error in predicting the ground truth top alternative, and error in predicting the ground truth ranking.

### 5.1 Difficulty & Expressiveness

We measure the following three metrics.

- *Response time*: Response time is known to be a good objective proxy for the cognitive load associated with a task [Rauterberg, 1992]. We measure the amount of time spent by the turkers on the tutorials and questions of the elicitation format.
- *Perceived difficulty*: As a subjective indicator of difficulty, we consider the perceived difficulty reported by the turkers (from Very Easy to Very Difficult) during the review stage of the elicitation format.
- *Perceived expressiveness*: Expressiveness indicates the amount of information that the turkers felt they were able to convey through the elicitation format (from Very Significant to Very Little).

Figure 2 shows the average time spent by the workers on the tutorial and on an average question under the six elicitation formats along with 95% confidence intervals (lower is better). We observe a statistically significant trend: when we fix a vote format (say Top or Rank), the average time spent increases for both tutorials and questions as we make the prediction format more complex (None  $\rightarrow$  Top  $\rightarrow$  Rank). In the full version, we show the average time spent for each domain and observe that the choice of the domain does not significantly affect it regardless of the elicitation format.

Figure 3 and Figure 4 respectively show the reported distributions of perceived difficulty (easier is better) and perceived expressiveness (higher is better). Interestingly, the turkers found the six elicitation formats to be of very similar difficulty and similar expressiveness.

### 5.2 Predicting the Ground Truth Top Alternative

We now turn to analyzing how effectively the different elicitation formats help us predict the ground truth. In addition to measuring the error of the ground truth estimate returned by our algorithm, we also measure the error in the input votes and predictions themselves. Note that every vote and prediction is an estimate of some truth (either the ground truth or a summary statistic of the other votes); thus, its error can be measured with respect to the truth it is attempting to uncover.

First, we consider predicting simply the top alternative in the ground truth ranking. For our algorithm as well as for the input votes and predictions, we use, as error measure, the frequency of incorrectly guessing the top alternative of the truth they attempt to estimate. Figure 5 shows the average prediction errors for various elicitation formats (lower is better).<sup>5</sup> We remind the reader that the effectiveness of SP voting should be judged based only on elicitation formats which include some prediction information.

Given four alternatives, selecting an alternative uniformly at random would result in a prediction error of 0.75. Interestingly, both the vote and prediction reports individually have average error around this benchmark. Yet, by combining these two pieces of individually erroneous information, SP voting is able to achieve significantly lower error. This is not surprising because SP voting approach is design precisely to pick out the minority of experts lurking among a majority of non-experts by combining vote and prediction information. Moreover, for a fixed type of vote (either Top or Rank), as the prediction formats become more complex (None  $\rightarrow$  Top  $\rightarrow$  Rank), the performance of SP voting improves.

Figure 6 compares SP voting to several standard voting rules including Plurality, Plurality with Runoff, Borda, Copeland, Instant Runoff Voting (IRV), and Maximin Rule, which ignore the prediction information and simply aggregate the vote information in a democratic manner.<sup>6</sup> The conventional voting rules run on elections containing votes from three elicitation formats (Rank-None, Rank-Top, and Rank-Rank) whereas SP voting runs on each elicitation format individually. We can see that for Rank-Rank, SP voting (rightmost orange bar) outperforms all conventional voting rules, despite having access to just a third of the samples. This indicates that the prediction information helps significantly.

These observations hold even when we consider each domain separately. These results are provided in the full version.

### 5.3 Predicting the Ground Truth Ranking

We now consider predicting the full ground truth ranking. For SP voting result as well as the individual votes and predictions, we use the Kendall-Tau (KT) distance to measure the error of the SP voting result, votes, and predictions compared to the true ranking they aim to estimate. Figure 7 shows the average KT distance for different elicitation formats (lower

<sup>5</sup>SP voting errors are obtained by averaging over 60 elections associated with 60 questions. Vote/Prediction errors are averaged over 1200 responses and have narrower confidence intervals.

<sup>6</sup>See [Brandt *et al.*, 2016] for definitions of these rules.

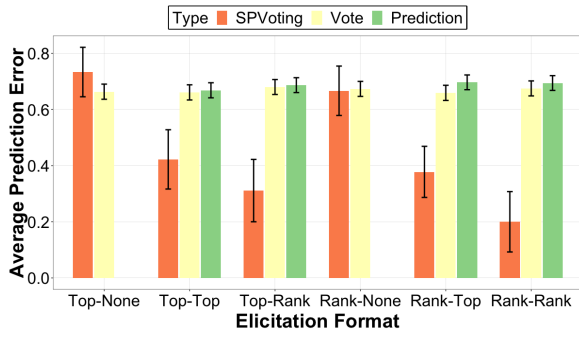


Figure 5: Average error in predicting the top alternative in the ground truth. By combining both the vote and predictions, SP voting achieves a much lower error than in either piece of information.

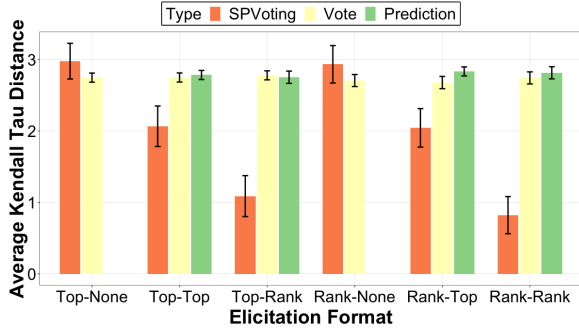


Figure 7: Average error in predicting the ground truth ranking. By combining both the vote and prediction information, SP voting achieves a much lower error than in either piece of information.

is better). Given four alternatives, selecting a uniformly random ranking will have an average KT distance of 3. Both the votes and prediction reports have average error around this benchmark. Similar to predicting the top alternative, SP voting produces significantly lower average error by combining these two noisy pieces of information. Moreover, for each vote format (either Top or Rank), as the prediction report becomes more expressive (None  $\rightarrow$  Top  $\rightarrow$  Rank) the average error of SP voting decreases.

Finally, we compare SP voting with standard voting rules (Figure 8) in terms of the average KT distance and find that SP voting again outperforms all voting rules for Rank-Rank.

### 5.4 Prediction vs. Vote

Our results illustrate the importance of prediction in recovering the ground truth. While eliciting ranked votes and predictions (Rank-Rank) achieves the lowest error, an intriguing question arises when we seek to choose an elicitation format that provides a reasonable tradeoff between accuracy and difficulty/expressiveness. Figures 5 and 7 show that Top-Rank significantly outperforms Rank-Top while both formats are comparable in terms of response time, perceived difficulty, and perceived expressiveness. Thus, if we wish to choose an elicitation format slightly more complex than Top-Top, making the prediction more expressive is more promising than that of the vote. The same observation holds when comparing Top-Top versus Rank-None. This shows that when a tradeoff

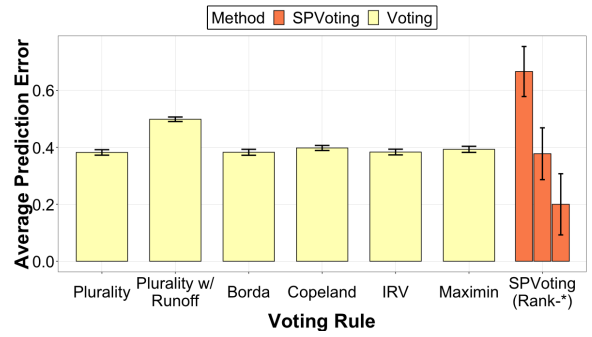


Figure 6: Comparing SP voting with conventional voting for predicting the top alternative. Incorporating the prediction reports helps SP voting significantly outperform conventional voting.

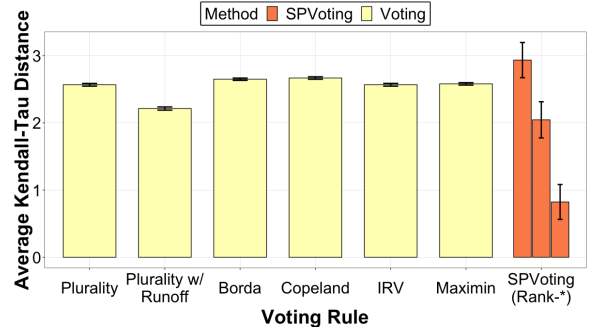


Figure 8: Comparing SP voting with conventional voting for predicting the ground truth ranking. Incorporating the prediction reports helps SP voting significantly outperform conventional voting.

between more complex vote and more complex prediction is necessary, eliciting more complex prediction may be better.

## 6 Discussion

We extended surprisingly popular voting to recover a ground truth ranking of alternatives and, through a crowdsourcing study across different domains, showed that it outperforms conventional voting approaches without significantly increasing elicitation. In our study, the ground truth is a ranking over four alternatives, and a challenging future direction is to extend this approach to rankings with more than four alternatives. For a large number of alternatives, any practical elicitation scheme would ask the voters to report a partial rank over the alternatives, which will make it challenging to design aggregation rules for such partial ranks.

Another interesting direction would be to derive theoretical performance guarantees for surprisingly popular voting when the number of participants is finite (the results of Prelec *et al.* [2017] hold only in the limit) and when only partial votes and predictions are elicited (this may require assuming a parametric signal distribution such as the Mallows model).

## Acknowledgements

The authors were partly supported by NSF grant #1850076 (Hosseini), a postdoctoral fellowship from Columbia DSI (Mandal), and an NSERC Discovery Grant (Shah).

## References

- [Abramowitz *et al.*, 2019] Ben Abramowitz, Elliot Anshelevich, and Wennan Zhu. Awareness of voter passion greatly improves the distortion of metric social choice. In *Proceedings of the 15th International Conference on Web and Internet Economics (WINE)*, pages 3–16, 2019.
- [Amanatidis *et al.*, 2020] Georgios Amanatidis, Georgios Birmpas, Aris Filos-Ratsikas, and Alexandros A. Voudouris. Peeking behind the ordinal curtain: Improving distortion via cardinal queries. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence (AAAI)*, pages 1782–1789, 2020.
- [Arrow *et al.*, 2008] Kenneth J. Arrow, Robert Forsythe, Michael Gorham, Robert Hahn, Robin Hanson, John O. Ledyard, Saul Levmore, Robert Litan, Paul Milgrom, Forrest D. Nelson, George R. Neumann, Marco Ottaviani, Thomas C. Schelling, Robert J. Shiller, Vernon L. Smith, Erik Snowberg, Cass R. Sunstein, Paul C. Tetlock, Philip E. Tetlock, Hal R. Varian, Justin Wolfers, and Eric Zitzewitz. The promise of prediction markets. *Science*, 320(5878):877–878, 2008.
- [Brandt *et al.*, 2016] Felix Brandt, Vincent Conitzer, Ulle Endriss, Jérôme Lang, and Ariel D Procaccia. *Handbook of computational social choice*. Cambridge University Press, 2016.
- [Camerer, 2011] Colin F Camerer. *Behavioral game theory: Experiments in strategic interaction*. Princeton University Press, 2011.
- [Caragiannis *et al.*, 2016] Ioannis Caragiannis, Ariel D Procaccia, and Nisarg Shah. When do noisy votes reveal the truth? *ACM Transactions on Economics and Computation (TEAC)*, 4(3):1–30, 2016.
- [Chen and Pennock, 2010] Yiling Chen and David M Pennock. Designing markets for prediction. *AI Magazine*, 31(4):42–52, 2010.
- [De Boer and Bernstein, 2017] Patrick M De Boer and Abraham Bernstein. Efficiently identifying a well-performing crowd process for a given problem. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing*, pages 1688–1699, 2017.
- [de Condorcet, 1785] Marquis de Condorcet. *Essai sur l’application de l’analyse à la probabilité de décisions rendues à la pluralité de voix*. Imprimerie Royal, 1785. Facsimile published in 1972 by Chelsea Publishing Company, New York.
- [Galesic *et al.*, 2018] Mirta Galesic, W Bruine de Bruin, Marion Dumas, A Kapteyn, JE Darling, and E Meijer. Asking about social circles improves election predictions. *Nature Human Behaviour*, 2(3):187–193, 2018.
- [Galton, 1907] Francis Galton. *Vox populi*. *Nature*, 75:450–451, 1907.
- [Kempe, 2020] David Kempe. Communication, distortion, and randomness in metric voting. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence (AAAI)*, pages 2087–2094, 2020.
- [Lalley and Weyl, 2018] Steven P Lalley and E Glen Weyl. Quadratic voting: How mechanism design can radicalize democracy. In *AEA Papers and Proceedings*, volume 108, pages 33–37, 2018.
- [Lee *et al.*, 2018] Michael D Lee, Irina Danileiko, and Julie Vi. Testing the ability of the surprisingly popular method to predict nfl games. *Judgment and Decision Making*, 13(4):322, 2018.
- [Mandal *et al.*, 2019] Debmalya Mandal, Ariel D Procaccia, Nisarg Shah, and David Woodruff. Efficient and thrifty voting by any means necessary. In *Advances in Neural Information Processing Systems*, pages 7180–7191, 2019.
- [Mandal *et al.*, 2020a] Debmalya Mandal, Goran Radanović, and David Parkes. The effectiveness of peer prediction in long-term forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 2160–2167, 2020.
- [Mandal *et al.*, 2020b] Debmalya Mandal, Nisarg Shah, and David P Woodruff. Optimal communication-distortion tradeoff in voting. In *Proceedings of the 21st ACM Conference on Economics and Computation*, pages 795–813, 2020.
- [Miller *et al.*, 2005] Nolan Miller, Paul Resnick, and Richard Zeckhauser. Eliciting informative feedback: The peer-prediction method. *Management Science*, 51(9):1359–1373, 2005.
- [Palley and Soll, 2019] Asa B Palley and Jack B Soll. Extracting the wisdom of crowds when information is shared. *Management Science*, 65(5):2291–2309, 2019.
- [Pivato, 2019] Marcus Pivato. Realizing epistemic democracy. In *The Future of Economic Design*, pages 103–112, 2019.
- [Prelec *et al.*, 2017] Dražen Prelec, H Sebastian Seung, and John McCoy. A solution to the single-question crowd wisdom problem. *Nature*, 541(7638):532, 2017.
- [Prelec, 2004] Dražen Prelec. A bayesian truth serum for subjective data. *science*, 306(5695):462–466, 2004.
- [Rauterberg, 1992] Matthias Rauterberg. A method of a quantitative measurement of cognitive complexity. *Human-computer interaction: Tasks and organisation*, pages 295–307, 1992.
- [Rutchick *et al.*, 2020] Abraham M Rutchick, Bryan J Ross, Dustin P Calvillo, and Catherine C Mesick. Does the “surprisingly popular” method yield accurate crowdsourced predictions? *Cognitive research: principles and implications*, 5(1):1–10, 2020.
- [Wang *et al.*, 2019] Juntao Wang, Yang Liu, and Yiling Chen. Forecast aggregation via peer prediction. arXiv:1910.03779, 2019.
- [Young, 1988] H. P. Young. Condorcet’s theory of voting. *The American Political Science Review*, 82(4):1231–1244, 1988.