

# On Optimality of Jury Selection in Crowdsourcing

Yudian Zheng, Reynold Cheng, Silviu Maniu, and Luyi Mo

Department of Computer Science, University of Hong Kong, Pokfulam Road, Hong Kong  
{ydzheng2, ckcheng, smaniu, lymo}@cs.hku.hk

## ABSTRACT

Recent advances in crowdsourcing technologies enable computationally challenging tasks (e.g., sentiment analysis and entity resolution) to be performed by Internet workers, driven mainly by monetary incentives. A fundamental question is: how should workers be selected, so that the tasks in hand can be accomplished successfully and economically? In this paper, we study the *Jury Selection Problem* (JSP): Given a monetary budget, and a set of decision-making tasks (e.g., “Is Bill Gates still the CEO of Microsoft now?”), return the set of workers (called *jury*), such that their answers yield the highest “Jury Quality” (or JQ). Existing JSP solutions make use of the *Majority Voting* (MV) strategy, which uses the answer chosen by the largest number of workers. We show that MV does not yield the best solution for JSP. We further prove that among all voting strategies (including deterministic and randomized strategies), *Bayesian Voting* (BV) can optimally solve JSP. We then examine how to solve JSP based on BV. This is technically challenging, since computing the JQ with BV is NP-hard. We solve this problem by proposing an approximate algorithm that is computationally efficient. Our approximate JQ computation algorithm is also highly accurate, and its error is proved to be bounded within 1%. We extend our solution by considering the task owner’s “belief” (or *prior*) on the answers of the tasks. Experiments on synthetic and real datasets show that our new approach is consistently better than the best JSP solution known.

## 1. INTRODUCTION

Due to advances in crowdsourcing technologies, computationally challenging tasks (e.g., sentiment analysis, entity resolution, document translation, etc.) can now be easily performed by human workers on the Internet. As reported by the Amazon Mechanical Turk in August 2012, over 500,000 workers from 190 countries worked on human intelligence tasks (HITs). The large number of workers and HITs have motivated researchers to develop solutions to streamline the crowdsourcing process [6,7,14,25,27,31,43,44].

In general, crowdsourcing a set of tasks involves the following steps: (1) distributing tasks to workers; (2) collecting the workers’ answers; (3) deciding final result; and (4) rewarding the workers.

An important question is: how should workers be chosen, so that the tasks in hand can be completed with high quality, while minimizing the monetary budget available? A related question, called the *Jury Selection Problem* (or JSP), has been recently proposed by Cao et al. [7]. Similar to the concept from law courts, a *jury*, or *jury set* denotes a subset of workers chosen from the available worker pool. Given a monetary budget and a task, the goal of JSP is to find the jury with the highest expected performance within the budget constraint. The kind of tasks studied in [7] is called the *decision-making task*: a question that requires an answer of either *yes* or *no* (e.g., “Is Bill Gates still the CEO of Microsoft now?”) and has a definitive ground truth. Decision-making tasks [7,39,41,44] are commonly used in crowdsourcing systems because of their conceptual simplicity. The authors of [7] were the first to propose a system to address JSP for this kind of tasks.

In this paper, we go beyond [7] and perform a comprehensive investigation of this problem. Particularly, we ask the following questions: (1) Is the solution in [7] optimal? (2) If not, what is an optimal solution for JSP? To understand these issues, let us first illustrate how [7] solves JSP.

Figure 1 shows a decision-making task, to be answered by some of the seven workers labeled from *A* to *G* where each worker is associated with a *quality* and a *cost*. The *quality* ranges from 0 to 1, indicating the probability that the worker correctly answers a question. This probability can be estimated by using her background information (e.g., her performance in other tasks) [7,25,37]. The *cost* is the amount of monetary reward the worker can get upon finishing a task. In this example, *A* has a quality of 0.77 and a cost of 9 units. For a jury, the *jury cost* is defined as the sum of workers’ costs in the jury and the *jury quality* (or JQ) is defined as the probability that the result returned by aggregating the jury answers is correct. Given a budget of *B* units, a feasible jury is a jury whose *jury cost* does not exceed *B*. For example, if  $B = \$20$ , then  $\{B, E, F\}$  is a feasible jury, since its *jury cost*, or  $\$5 + \$5 + \$2 = \$12$ , is not larger than  $\$20$ .

To solve JSP, a naive solution is to compute the JQ for every feasible jury, and return the one with the highest JQ. [7] studies how to compute JQ for a jury where the jury’s returned result is decided by *Majority Voting* (MV). In short, MV returns the result as the one corresponding to the most workers. In the following, we consider each worker’s answer as a “vote” for either “yes” or “no”. Let us consider  $\{B, E, F\}$  again, the probability that these workers gives a correct result according to MV is  $0.7 \cdot 0.6 \cdot 0.6 + 0.7 \cdot 0.6 \cdot (1 - 0.6) + 0.7 \cdot (1 - 0.6) \cdot 0.6 + (1 - 0.7) \cdot 0.6 \cdot 0.6 = 69.6\%$ . Moreover, since  $\{A, C, G\}$  yields the highest JQ among all the feasible juries, it is considered to be the best solution by [7].

As illustrated above, MV is used to solve JSP in [7]. In addition to MV, researchers have proposed a variety of *voting strategies*,

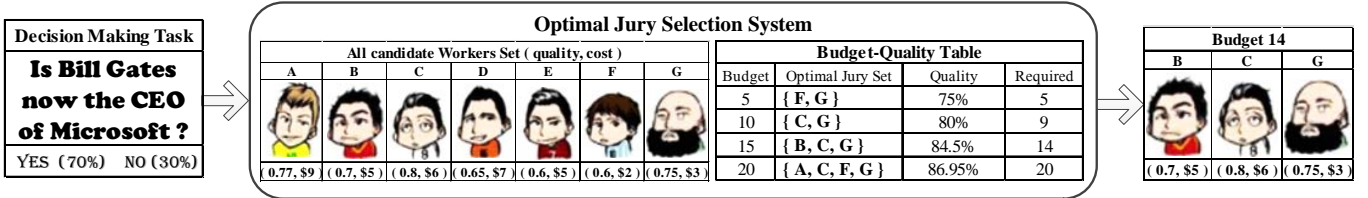


Figure 1: Optimal Jury Selection System.

such as Bayesian Voting (BV) [25], Randomized Majority Voting [20], and Random Ballot Voting [33]. Like MV, these voting strategies decide the final result of a decision-making task based on the workers’ votes. For example, BV computes the posterior probability of answers according to Bayes’ Theorem [3], based on the workers’ votes, and returns the answer having the largest posterior probability.

In this paper, we investigate an interesting problem: is it possible to find the optimal voting strategy for JSP among all voting strategies? One simple answer to this question is to consider all voting strategies. However, as listed in Table 2, the number of existing strategies is very large. Moreover, multiple new strategies may emerge in the future. We address this question by first studying the criteria of a strategy that produce an optimal solution for JSP (i.e., given a jury, the JQ of the strategy is the highest among all the possible voting strategies). This is done by observing that voting strategies can be classified into two major categories: *deterministic* and *randomized*. A deterministic strategy aggregates workers’ answers without any degree of randomness; MV is a typical example of this class. For a randomized strategy, each answer is returned with some probability. Using this classification, we present the criteria required for a voting strategy that leads to the optimal solution for JSP. We discover that BV satisfies the requirements of an optimal strategy. In other words, BV is the optimal voting strategy with respect to JQ, and will consistently produce better quality juries than the other strategies.

How to solve JSP with BV then? A straightforward solution is to enumerate all feasible juries, and find the one with the largest value of JQ. However, this approach suffers from two major problems:

1. Computing the JQ of a jury for BV requires enumerating an exponentially large number of workers’ answers. In fact, we show that this problem is NP-hard;
2. The number of feasible juries is exponentially large.

To solve Problem 1, we develop a polynomial-time-based approximation algorithm, which enables a large number of candidate answers to be pruned, without a significant loss of accuracy. We further develop a theoretical error bound of this algorithm. Particularly, our approximate JQ computation algorithm is proved to yield an error of not more than 1%. To tackle Problem 2, we leverage a successful heuristic, the simulated annealing heuristic, by designing local neighborhood search functions. To evaluate our solutions, we have performed extensive evaluation on real and synthetic crowdsourced data. Our experimental results show that our algorithms effectively and efficiently solve JSP. The quality of our solution is also consistently better than that of [7].

We also study how to allow the provider of the tasks to place her confidence information (called *prior*) on the answers of the task. She may associate a “belief score” on the answers to the tasks, before the crowdsourcing process starts. For instance, in Figure 1, if she is more confident that Bill Gates is still the CEO of Microsoft, she can assign 70% to *yes*, and 30% to *no*. Intuitively, we prove

that under BV, the effect of prior is just the same as regarding the task provider as another worker, having the same quality values as the prior.

Figure 1 illustrates our crowdsourcing system, which we called the “Optimal Jury Selection System”. In this system, the task provider published a decision-making task. Then, based on the workers’ information (i.e., their individual quality and cost), a “budget-quality table” is generated. In this table, each row contains a budget, the computed optimal jury, its estimated jury quality and the required budget for the jury. Based on this table, the task provider can conveniently decide the best budget-quality combination. For example, she may deem that increasing the budget from 15 units to 20 units is not worthwhile, since the quality increases only by around 2.5%. In this example, the task provider selects the jury set  $\{B, C, G\}$  that is the best under a budget of 15 units. This chosen jury set would cost her only 14 units.

Recall that [7] focuses on addressing JSP under MV on decision-making tasks and we address the optimality of JSP on decision-making tasks by considering all voting strategies, where each worker’s quality is modeled as a single parameter. In reality, multiple choice tasks [25,34,42] are also commonly used in crowdsourcing and several works [1,34,36] model each worker as a confusion matrix rather than a single quality score. We also briefly discuss here the optimality of JSP for other task types and worker models, and how our solutions can be extended to these other variants.

The rest of this paper is arranged as follows. We describe the data model and the problem definition in Section 2. In Section 3, we examine the requirements of an optimal voting strategy for JSP, and show that BV satisfies these criteria. We present an efficient algorithm to compute JQ of a jury set in Section 4 and develop fast solutions to solve JSP in Section 5. In Section 6, we present our experimental results. We discuss how our solutions can be extended for other task types and worker models in Section 7. In Section 8, we review the related works and Section 9 concludes the paper.

## 2. DATA MODEL & PROBLEM DEFINITION

We now describe our data model in Section 2.1 and define the jury selection problem in Section 2.2.

### 2.1 Data Model

In this paper, we focus on the *decision-making tasks* where each task has two possible answers (either *yes* or *no*). We use 1 and 0 to denote *yes* and *no*, respectively. We assume that each task has a latent true answer (or ground truth)  $t \in \{0, 1\}$ , which is unknown in advance. The task provider usually assigns a *prior* on the task, which describes her prior knowledge in the probability distribution of the task’s true answer. We denote the prior by  $\alpha$  where  $\Pr(t = 0) = \alpha$ , and  $\Pr(t = 1) = 1 - \alpha$ . If the task provider has no prior knowledge for the task, then we assume  $\alpha = 0.5$ .

A *jury* (or *jury set*), denoted by  $J$ , is a collection of  $n$  workers drawn from a set of  $N$  candidate workers  $W = \{j_1, j_2, \dots, j_N\}$ , i.e.,  $J \subseteq W$ ,  $|J| = n$ . Without loss of generality, let  $J =$

$\{j_1, j_2, \dots, j_n\}$ . In order to infer the ground truth ( $\mathbf{t}$ ), we leverage the collective intelligence of a jury, i.e. we ask each worker to give a vote for the task. We use  $V$ , a *voting*, to denote the set of votes (answers) given by a jury  $J$ , and so  $V = \{v_1, v_2, \dots, v_n\}$  where  $v_i \in \{0, 1\}$  is the vote given by  $j_i$ . We assume the independence of each worker’s vote, an assumption also used in [7,18,25,34].

We follow the worker model in previous works [7,25,44], where each worker  $j_i$  is associated with a quality  $q_i \in [0, 1]$  and a cost  $c_i$ . The quality  $q_i$  indicates the probability that the worker conducts a correct vote, i.e.,  $q_i = \Pr(v_i = \mathbf{t})$ , and the cost  $c_i$  represents the money (or incentive) required for  $j_i$  to give a vote. A few works [7,25,37] have recently addressed how to derive the quality and the cost of a worker by leveraging the backgrounds and answering history of individuals. Thus, similar to [7], we assume that they are known in advance.

We remark that the optimality of JSP and our solutions can be extended to address other task types and worker models used in [1,25,34,34,36,42]. We will briefly discuss these extensions in Section 7.

## 2.2 Problem Definition

Let  $B$  be the budget of a task provider, i.e., a maximum of  $B$  cost units can be given to a jury to collect their votes. Our goal is to solve the *Jury Selection Problem* (denoted by JSP) which selects a jury  $J$  under the budget constraint ( $\sum_{j_i \in J} c_i \leq B$ ) such that the jury’s collective intelligence is maximized.

The collective intelligence of a jury is closely related to the *Voting Strategy*, denoted by  $S$ , which estimates the true answer of the task based on the prior, the jury and their votes. We say the estimated true answer is the *result* of the voting strategy. A detailed discussion about the voting strategy is given in Section 3.1.

In order to quantify the jury’s collective intelligence, we define the *Jury Quality* (or  $JQ$  in short) which essentially measures the probability that the result of the voting strategy is correct. The score of  $JQ$  is given by function  $JQ(J, S, \alpha)$ . We will give a precise definition for  $JQ$  in Section 3.2.

Let  $\Theta$  denote the set of all voting strategies and  $\mathcal{C}$  denote the set of all feasible juries (i.e.,  $\mathcal{C} = \{J \mid J \subseteq W \wedge \sum_{j_i \in J} c_i \leq B\}$ ). The aim of JSP is to select the optimal jury  $J^*$  such that

$$\text{given } \alpha \text{ and } q_i, c_i \text{ (for } i = 1, 2, \dots, N) \quad (1)$$

$$J^* = \arg \max_{J \in \mathcal{C}} \max_{S \in \Theta} JQ(J, S, \alpha) \quad (2)$$

Note that existing work [7] only focuses on majority voting strategy (MV) and solves  $\arg \max_{J \in \mathcal{C}} JQ(J, MV, 0.5)$ , which, as we shall prove later, is sub-optimal for JSP.

In the rest of the paper, we first discuss how to derive the optimal voting strategy  $S^*$  such that  $JQ(J, S^*, \alpha) = \max_{S \in \Theta} JQ(J, S, \alpha)$  (Section 3). We then talk about the computation of  $JQ(J, S^*, \alpha)$  (Section 4) and finally address the problem of finding  $J^*$  (Section 5).

Table 1 summarizes the symbols used in this paper.

## 3. OPTIMAL VOTING STRATEGY

In this section, we present a detailed description for the voting strategy in Section 3.1. We then formally define  $JQ$  in Section 3.2. Finally, we give an optimal voting strategy with respect to  $JQ$  in Section 3.3.

### 3.1 Voting Strategies

As mentioned, a *voting strategy*  $S$  gives an estimation of the true answer  $\mathbf{t}$  based on the prior  $\alpha$ , the jury  $J$  and their votes  $V$ . Thus, we model a voting strategy as a function  $S(V, J, \alpha)$ , whose result is an estimation of  $\mathbf{t}$ . Based on whether the result is given with

**Table 1: Table of Symbols**

Symbol	Description
$\mathbf{t}$	the ground truth for a task, and $\mathbf{t} \in \{0, 1\}$
$\alpha$	prior given by the task provider, and $\alpha = \Pr(\mathbf{t} = 0)$
$W$	a set of all candidate workers $W = \{j_1, j_2, \dots, j_N\}$
$J$	a jury, $J \subseteq W$ and $ J  = n$ , $J = \{j_1, j_2, \dots, j_n\}$
$V$	a voting given by $J$ , and $V = \{v_1, v_2, \dots, v_n\}$
$q_i$	the quality of worker $j_i$ and $q_i \in [0, 1]$
$c_i$	the cost of worker $j_i$
$B$	the budget provided by the task provider
$\Theta$	a set containing all voting strategies
$\mathcal{C}$	the set of all possible juries within budget constraint

degree of randomness, we can classify the voting strategies into two categories: *deterministic voting strategy* and *randomized voting strategy*.

**DEFINITION 1.** A *deterministic voting strategy*  $S(V, J, \alpha)$  returns the result as 0 or 1 without any degree of randomness.

**DEFINITION 2.** A *randomized voting strategy*  $S(V, J, \alpha)$  returns the result as 0 with probability  $p$  and 1 with probability  $1 - p$ .

**EXAMPLE 1.** The *majority voting strategy* (or *MV*) is a typical *deterministic voting strategy*, and it gives result as 0 if more than half of workers vote for 0 (i.e.,  $\sum_{i=1}^n (1 - v_i) \geq \frac{n+1}{2}$ ); otherwise, the result is 1.

Its *randomized counterpart* is called *randomized majority voting strategy* (or *RMV*), which returns the result with probability proportional to the number of votes. That is, *RMV* returns 0 with probability  $p = \frac{1}{n} \sum_{i=1}^n (1 - v_i)$ , and 1 with probability  $1 - p$ .

Note that randomized strategies are often introduced to improve the error bound for worst-case analysis [23]. And thus, they are widely used when the worst-case performance is the main concern.

**Table 2: Classification of Voting Strategies**

Deterministic Voting Strategies	Randomized Voting Strategies
Majority Voting (MV) [7]	Randomized Majority Voting (RMV) [20]
Half Voting [28]	Random Ballot Voting [33]
Bayesian Voting [25]	Triadic Consensus [2]
Weighted MV [23]	Randomized Weighted MV [23]
...	...

Table 2 shows a few voting strategies, which are introduced in previous works, and their corresponding category.

### 3.2 Jury Quality

In order to measure the goodness of a voting strategy  $S$  for a jury  $J$ , we introduce a metric called *Jury Quality* (or  $JQ$  in short). We model  $JQ$  by a function  $JQ(J, S, \alpha)$  which gives the quality score as the probability of drawing a correct result under the voting strategy, i.e.,

$$JQ(J, S, \alpha) = \Pr(S(\mathbf{V}, J, \alpha) = \mathbf{t}) \quad (3)$$

where  $\mathbf{V} \in \{0, 1\}^n$  and  $\mathbf{t} \in \{0, 1\}$  are two random variables corresponding to the unknown jury’s voting, and the task’s latent true answer. For notational convenience, we omit  $J$  and  $\alpha$  in  $S$  when their values are understood and simply write  $S(\mathbf{V})$  instead of  $S(\mathbf{V}, J, \alpha)$ .

Let  $\mathbb{1}_{\{st\}}$  be the indicator function, which returns 1 if the statement  $st$  is true, and 0 otherwise. Let  $\Omega$  be the domain of  $\mathbf{V}$ , i.e.,

$\Omega = \{0, 1\}^n$ .  $JQ(J, S, \alpha)$  can be rewritten as follows.

$$\begin{aligned} JQ(J, S, \alpha) &= 1 \cdot \Pr(S(\mathbf{V}) = \mathbf{t}) + 0 \cdot \Pr(S(\mathbf{V}) \neq \mathbf{t}) \\ &= \mathbb{E}[\mathbb{1}_{\{S(\mathbf{V})=\mathbf{t}\}}] \\ &= \sum_{\mathbf{t} \in \{0,1\}} \sum_{V \in \Omega} \Pr(\mathbf{V} = V, \mathbf{t} = t) \cdot \mathbb{E}[\mathbb{1}_{\{S(V)=t\}}] \end{aligned}$$

We now give a precise definition for JQ as below.

**DEFINITION 3 (JURY QUALITY).** *Given a jury  $J$  and the prior  $\alpha$ , the Jury Quality (or JQ) for a voting strategy  $S$ , denoted by  $JQ(J, S, \alpha)$ , is defined as*

$$\begin{aligned} &\alpha \cdot \sum_{V \in \Omega} \Pr(\mathbf{V} = V | \mathbf{t} = 0) \cdot \mathbb{E}[\mathbb{1}_{\{S(V)=0\}}] \\ &+ (1 - \alpha) \cdot \sum_{V \in \Omega} \Pr(\mathbf{V} = V | \mathbf{t} = 1) \cdot \mathbb{E}[\mathbb{1}_{\{S(V)=1\}}]. \end{aligned} \quad (4)$$

For notational convenience, we write  $\Pr(V|\mathbf{t} = 0)$  instead of  $\Pr(\mathbf{V} = V|\mathbf{t} = 0)$ , and  $\Pr(V|\mathbf{t} = 1)$  instead of  $\Pr(\mathbf{V} = V|\mathbf{t} = 1)$ . Next, we give two marks in computing JQ.

1. Since workers give votes independently, we have

$$\begin{aligned} \Pr(V | \mathbf{t} = 0) &= \prod_{i=1}^n q_i^{(1-v_i)} \cdot (1 - q_i)^{v_i} \\ \Pr(V | \mathbf{t} = 1) &= \prod_{i=1}^n q_i^{v_i} \cdot (1 - q_i)^{(1-v_i)} \end{aligned}$$

2.  $\mathbb{E}[\mathbb{1}_{\{S(V)=0\}}]$  and  $\mathbb{E}[\mathbb{1}_{\{S(V)=1\}}]$  are either 0 or 1 if  $S$  is a deterministic voting strategy; or value of  $p$  and  $1 - p$  if  $S$  is a randomized voting strategy (refer to Definition 2).

We next give an example to illustrate the computation of JQ.

**EXAMPLE 2.** *Suppose  $\alpha = 0.5$  and there are 3 workers in  $J$  with workers' qualities as 0.9, 0.6, 0.6 respectively. To compute JQ for MV, we enumerate all possible combinations of  $V \in \{0, 1\}^3$  and  $t \in \{0, 1\}$ , and show the results in Figure 2. The 3rd column in each table represents the probability that a specific combination ( $V$  and  $t$ ) exists. The 4th column shows the result of MV for each  $V$ . The symbol  $\checkmark$  indicates whether MV's result is correct or not (according to the value of  $t$ ). And thus,  $JQ(J, MV, \alpha)$  equals to the summation of probabilities where symbol  $\checkmark$  occurs. Take  $V = \{1, 0, 0\}$  and  $t = 0$  as an example. First,  $\Pr(\mathbf{V} = V, \mathbf{t} = 0) = 0.018$ . Since  $\sum_{i=1}^3 (1 - v_i) = 2 \geq \frac{n+1}{2} = 2$ , we have  $MV(V) = 0 = t$ . Thus, the probability 0.018 is added to  $JQ(J, MV, \alpha)$ . Similarly, for  $V = \{1, 0, 0\}$  and  $t = 1$ , as  $MV(V) = 0 \neq t$ , then  $\Pr(\mathbf{V} = V, \mathbf{t} = 1) = 0.072$  will not be added to  $JQ(J, MV, \alpha)$ . Considering all  $V$ 's and  $t$ 's, the final  $JQ(J, MV, \alpha) = 79.2\%$ .*

### 3.3 Optimal Voting Strategy

In the last two sections, we present a few voting strategies and define Jury Quality to quantify the goodness of a voting strategy. Thus an interesting question is: does an optimal voting strategy  $S^*$  with respect to JQ exist? That is, given any  $J$  and  $\alpha$ ,  $JQ(J, S^*, \alpha) = \max_{S \in \Theta} JQ(J, S, \alpha)$ . Note that if  $S^*$  exists, we can then solve JSP without enumerating all voting strategies in  $\Theta$  (refer to Equation 2).

To answer this question, let us reconsider Definition 3. Let  $h(V) = \mathbb{E}[\mathbb{1}_{\{S(V)=0\}}]$ . We have (i)  $h(V) \in [0, 1]$ ; and (ii)  $\mathbb{E}[\mathbb{1}_{\{S(V)=1\}}] = 1 - h(V)$ . Also, let  $P_0(V) = \Pr(\mathbf{V} = V, \mathbf{t} = 0)$ , and  $P_1(V) = \Pr(\mathbf{V} = V, \mathbf{t} = 1)$ . Hence,  $JQ(J, S, \alpha)$  can be rewritten as

$$\begin{aligned} &\sum_{V \in \Omega} [P_0(V) \cdot h(V) + P_1(V) \cdot (1 - h(V))] \\ &= \sum_{V \in \Omega} [h(V) \cdot (P_0(V) - P_1(V)) + P_1(V)] \end{aligned}$$

No.	$V$	$\Pr(\mathbf{t}=0) \cdot \Pr(\mathbf{V}=V   \mathbf{t}=0)$	MV: [compare] (result) [✓/✗]	BV: [compare] (result) [✓/✗]
1	{0,0,0}	0.5*0.9*0.6*0.6=0.162	[3 ≥ 2] (0) [✓]	[0.162 ≥ 0.008] (0) [✓]
2	{0,0,1}	0.5*0.9*0.6*0.4=0.108	[2 ≥ 2] (0) [✓]	[0.108 ≥ 0.012] (0) [✓]
3	{0,1,0}	0.5*0.9*0.4*0.6=0.108	[2 ≥ 2] (0) [✓]	[0.108 ≥ 0.012] (0) [✓]
4	{0,1,1}	0.5*0.9*0.4*0.4=0.072	[1 < 2] (1) [✗]	[0.072 ≥ 0.018] (0) [✓]
5	{1,0,0}	0.5*0.1*0.6*0.6=0.018	[2 ≥ 2] (0) [✓]	[0.018 < 0.072] (1) [✗]
6	{1,0,1}	0.5*0.1*0.6*0.4=0.012	[1 < 2] (1) [✗]	[0.012 < 0.018] (1) [✗]
7	{1,1,0}	0.5*0.1*0.4*0.6=0.012	[1 < 2] (1) [✗]	[0.012 < 0.018] (1) [✗]
8	{1,1,1}	0.5*0.1*0.4*0.4=0.008	[0 < 2] (1) [✗]	[0.008 < 0.162] (1) [✗]

(a) Enumeration of all  $2^3 = 8$  possible votings in  $\Omega$  ( $\mathbf{t} = 0$ )

No.	$V$	$\Pr(\mathbf{t}=1) \cdot \Pr(\mathbf{V}=V   \mathbf{t}=1)$	MV: [compare] (result) [✓/✗]	BV: [compare] (result) [✓/✗]
1	{0,0,0}	0.5*0.1*0.4*0.4=0.008	[3 ≥ 2] (0) [✗]	[0.162 ≥ 0.008] (0) [✗]
2	{0,0,1}	0.5*0.1*0.4*0.6=0.012	[2 ≥ 2] (0) [✗]	[0.108 ≥ 0.012] (0) [✗]
3	{0,1,0}	0.5*0.1*0.6*0.4=0.012	[2 ≥ 2] (0) [✗]	[0.108 ≥ 0.012] (0) [✗]
4	{0,1,1}	0.5*0.1*0.6*0.6=0.018	[1 < 2] (1) [✓]	[0.072 ≥ 0.018] (0) [✗]
5	{1,0,0}	0.5*0.9*0.4*0.4=0.072	[2 ≥ 2] (0) [✗]	[0.018 < 0.072] (1) [✓]
6	{1,0,1}	0.5*0.9*0.4*0.6=0.108	[1 < 2] (1) [✓]	[0.012 < 0.018] (1) [✓]
7	{1,1,0}	0.5*0.9*0.6*0.4=0.108	[1 < 2] (1) [✓]	[0.012 < 0.018] (1) [✓]
8	{1,1,1}	0.5*0.9*0.6*0.6=0.162	[0 < 2] (1) [✓]	[0.008 < 0.162] (1) [✓]

(b) Enumeration of all  $2^3 = 8$  possible votings in  $\Omega$  ( $\mathbf{t} = 1$ )

**Figure 2: Example of JQ computation for MV and BV ( $\alpha = 0.5$ , and the quality of workers are 0.9, 0.6, and 0.6)**

This gives us a hint to maximize  $JQ(J, S, \alpha)$  and find the optimal voting strategy  $S^*$ . Let  $h^*(V) = \mathbb{E}[\mathbb{1}_{\{S^*(V)=0\}}]$ . It is observed that  $P_1(V)$  is constant for a given  $V$  and  $h(V) \in [0, 1]$  for all  $S$ 's (no matter it is a deterministic one or a randomized one). Thus, to optimize  $JQ(J, S, \alpha)$ , it is required that

1. if  $P_0(V) - P_1(V) < 0$ ,  $h^*(V) = 0$ , and so,  $S^*(V) = 1$ ;
2. if  $P_0(V) - P_1(V) \geq 0$ ,  $h^*(V) = 1$ , and so,  $S^*(V) = 0$ .

We summarize this observation as below.

**THEOREM 1.** *Given  $\alpha$ ,  $J$ , and  $V$ , the optimal voting strategy, denoted by  $S^*$ , decides the result as follows:*

1.  $S^*(V) = 1$  if  $\alpha \cdot \prod_{i=1}^n q_i^{(1-v_i)} \cdot (1 - q_i)^{v_i} < (1 - \alpha) \cdot \prod_{i=1}^n q_i^{v_i} \cdot (1 - q_i)^{(1-v_i)}$ ; or
2.  $S^*(V) = 0$ , otherwise.

Note that  $S^*$  is a deterministic voting strategy, and it's essentially a voting strategy based on the Bayes' Theorem [11]. The reason is as follows. According to the Bayes' Theorem, based on the observed voting  $V$ ,  $\Pr(\mathbf{t} = 0 | \mathbf{V} = V) = P_0(V) / \Pr(\mathbf{V} = V)$ , and similarly  $\Pr(\mathbf{t} = 1 | \mathbf{V} = V) = P_1(V) / \Pr(\mathbf{V} = V)$ . Therefore,  $P_0(V) - P_1(V) < 0$  indicates  $\Pr(\mathbf{t} = 0 | \mathbf{V} = V) < \Pr(\mathbf{t} = 1 | \mathbf{V} = V)$ . And so, 1 has a higher probability to be the true answer than 0. Thus, the voting strategy based on the Bayes' Theorem returns 1 as the result, which is consistent with  $S^*$  in Theorem 1. Next, we give a formal definition for Bayesian Voting (BV) and summarize the above observation in Theorem 1.

**DEFINITION 4.** *The voting strategy based on the Bayes' Theorem, denoted by Bayesian Voting (or BV in short), returns the result as 1, if  $\Pr(\mathbf{t} = 0) \cdot \Pr(\mathbf{V} = V | \mathbf{t} = 0) < \Pr(\mathbf{t} = 1) \cdot \Pr(\mathbf{V} = V | \mathbf{t} = 1)$ ; or 0, otherwise.*

**COROLLARY 1.** *BV is optimal w.r.t. JQ, i.e.,  $S^* = BV$ .*

Note that the BV is also used in [1, 18, 25]. In the rest of the paper, we use  $S^*$  and BV interchangeably. We remark that the optimality of BV is based on two assumptions: (1) the prior and workers' qualities are known in advance; (2) JQ (Definition 3) is adopted to measure the goodness of a voting strategy.

EXAMPLE 3. Let us reconsider Figure 2 and see how  $JQ(J, BV, \alpha)$  is computed. The 5th column shows results given by BV. The two numbers in bracket correspond to  $P_0(V)$  and  $P_1(V)$ , respectively. The value in parenthesis is the estimated true answer returned by BV. We again use a symbol  $\checkmark$  to indicate the correct voting result. Take  $V = \{1, 0, 0\}$  and  $t = 0$  as an example. Since  $\alpha \cdot (1 - q_1) \cdot q_2 \cdot q_3 = 0.018 < (1 - \alpha) \cdot q_1 \cdot (1 - q_2) \cdot (1 - q_3) = 0.072$ , we have  $BV(V) = 1 \neq t$ , thus 0.018 is not added into  $JQ(J, BV, \alpha)$ . Otherwise, for  $V = \{1, 0, 0\}$  and  $t = 0$ , similarly we derive that 0.072 is added in  $JQ(J, BV, \alpha)$ . Recall Example 2, when  $V = \{1, 0, 0\}$ , if we consider two cases of  $t$ , then 0.072 is added into  $JQ(J, BV, \alpha)$ ; but here we have seen in Example 2 that 0.018 is added into  $JQ(J, MV, \alpha)$ . By considering all  $V$  and  $t$ , we have  $JQ(J, BV, \alpha) = 90\% > JQ(J, MV, \alpha) = 79.2\%$ .

Intuitively, the reason why BV outperforms other voting strategies is that BV considers the prior and worker's qualities in deriving the result of a voting  $V$ , and only the one with larger posterior probability is returned. Thus, it is more likely to return a correct answer than other strategies. For example, assume  $\alpha = 0.5$  and the voting  $V = \{0, 1, 1\}$  is given by workers with individual quality 0.9, 0.6 and 0.6 respectively. As  $0.5 \cdot 0.9 \cdot (1 - 0.6) \cdot (1 - 0.6) > 0.5 \cdot (1 - 0.9) \cdot 0.6 \cdot 0.6$ , BV returns 0 as the result. However, MV does not leverage either the prior information or workers' qualities, and so, it returns 1, which is given by two lower quality workers.

Before we move on, we would like to discuss the effect of  $q_i$  for voting strategies. Intuitively,  $q_i < 0.5$  indicates that worker  $j_i$  is more likely to give an incorrect answer than a correct one. Thus, we can either simply ignore this worker in the jury selection process, or modify her answer according to the specific voting strategy. For example, for MV, we can regard vote 0 as 1 and vote 1 as 0 if the vote is given by a worker whose quality is less than 0.5; for BV, according to its definition, it can reinterpret the vote given by a worker with quality  $q_i < 0.5$  as an opposite vote given by a worker with quality  $1 - q_i > 0.5$ .<sup>1</sup> Moreover, in our experiments with real human workers, we observed that their qualities were generally well above 0.5. We thus assume that  $q_i \geq 0.5$  in our subsequent discussions, without loss of generality.

## 4. COMPUTING JURY QUALITY FOR OPTIMAL STRATEGY

In the previous section, we have proved that BV is the optimal voting strategy with respect to JQ. And thus, in order to solve JSP, we only need to figure out  $J^*$  such that  $JQ(J^*, BV, \alpha)$  is maximized. An immediate question is whether  $JQ(J, BV, \alpha)$  can be computed efficiently. Unfortunately, we find that computing  $JQ(J, BV, \alpha)$  is NP-hard (Section 4.1). To alleviate this, we propose an efficient approximation algorithm with theoretical bounds to compute JQ for BV in this section.

### 4.1 NP-hardness of computing $JQ(J, BV, \alpha)$

Note that [7] has previously proposed an efficient algorithm to compute  $JQ(J, MV, 0.5)$  in  $\mathcal{O}(n \log n)$ . However, this polynomial algorithm cannot be adapted to compute JQ for BV. The main reason is that computing JQ for BV is an NP-hard problem.

**THEOREM 2.** Given  $\alpha$  and  $J$ , computing JQ for BV, or  $JQ(J, BV, \alpha)$ , is NP-hard.

The proof is non-trivial and we present the detailed proof in the technical report [15] due to space limits. The idea of

<sup>1</sup>Details of the reinterpretation can be found in the technical report [15].

$R(V) = u(V) - w(V)$	$A_0(V)$	$A_1(\bar{V})$	$A_0(V) + A_1(\bar{V})$
$R(V) > 0$	$e^{u(V)}/2$	$e^{u(V)}/2$	$e^{u(V)}$
$R(V) = 0$	$e^{u(V)}/2$	0	$e^{u(V)}/2$
$R(V) < 0$	0	0	0

Figure 3: Expressing  $A_0(V) + A_1(\bar{V})$  using  $R(V)$  and  $u(V)$

the proof is that the partition problem [32] (a well-known NP-complete problem) can be reduced to the problem of computing  $JQ(J, BV, 0.5)$  for some  $J$ . Hence, the computation of  $JQ(J, BV, 0.5)$  is not easier than the partition problem. Moreover, computing  $JQ(J, BV, 0.5)$  is not in NP (it is not a decision problem), which makes the problem of computing  $JQ(J, BV, \alpha)$  for  $\alpha \in [0, 1]$  NP-hard.

To avoid this hardness result, we propose an approximation algorithm. We first discuss the computation of  $JQ(J, BV, 0.5)$  in Section 4.2 and 4.3, and give its approximation error bound in Section 4.4. Finally, we briefly discuss how to adapt the algorithm to  $\alpha \in [0, 1]$  in Section 4.5.

### 4.2 Analysis of Computing $JQ(J, BV, 0.5)$

Let us first give some basic analysis for computing  $JQ(J, BV, 0.5)$  before we introduce our approximation algorithm. To facilitate our analysis, we first define a few symbols.

- $A_0(V) = 0.5 \cdot \Pr(V | t = 0) \cdot \mathbb{1}_{\{BV(V)=0\}}$ ;
- $A_1(V) = 0.5 \cdot \Pr(V | t = 1) \cdot \mathbb{1}_{\{BV(V)=1\}}$ ;
- $\bar{V} = \{\bar{v}_1, \bar{v}_2, \dots, \bar{v}_n\}$ , where  $\bar{v}_i = 1 - v_i$  ( $1 \leq i \leq n$ ).

From Figure 2 we observe that  $A_0(V) = A_1(\bar{V})$ . For example,  $A_0(\{0, 1, 0\}) = A_1(\{1, 0, 1\}) = 0.108$  and  $A_0(\{1, 0, 1\}) = A_1(\{0, 1, 0\}) = 0$ . The observation motivates us to consider  $A_0(V)$  and  $A_1(\bar{V})$  together, and we can prove that

$$\begin{aligned} JQ(J, BV, 0.5) &= \sum_{V \in \Omega} [A_0(V) + A_1(V)] \\ &= \sum_{V \in \Omega} [A_0(V) + A_1(\bar{V})], \end{aligned} \quad (5)$$

as  $V \rightarrow \bar{V}$  defines a one-to-one correspondence between  $\Omega$  and  $\Omega$ .

We further define  $u(V)$  and  $w(V)$  as follows.

$$\begin{aligned} u(V) &= \ln \Pr(V | t = 0) = \sum_{i=1}^n [(1 - v_i) \ln q_i + v_i \ln(1 - q_i)], \\ w(V) &= \ln \Pr(V | t = 1) = \sum_{i=1}^n [v_i \ln q_i + (1 - v_i) \ln(1 - q_i)], \end{aligned}$$

Let  $R(V) = u(V) - w(V)$  and  $\sigma(q_i) = \ln \frac{q_i}{1 - q_i}$  (as  $q_i \geq 0.5$ ,  $\sigma(q_i) \geq 0$ ), we have

$$R(V) = \sum_{i=1}^n [(1 - 2v_i) \cdot \sigma(q_i)], \quad e^{u(V)} = \prod_{i=1}^n q_i^{(1 - v_i)} \cdot (1 - q_i)^{v_i}. \quad (6)$$

As illustrated in Figure 3, we can express  $A_0(V) + A_1(\bar{V})$  based on the sign of  $R(V)$  and the value of  $u(V)$ . And therefore,

$$JQ(J, BV, 0.5) = \sum_{V \in \Omega} [\mathbb{1}_{\{R(V) > 0\}} \cdot e^{u(V)} + \mathbb{1}_{\{R(V) = 0\}} \cdot \frac{e^{u(V)}}{2}].$$

Motivated by the above formula<sup>2</sup>, we can apply an iterative approach which expands  $J$  with one more worker at each iteration and thus compute  $JQ(J, BV, 0.5)$  in  $n$  total iterations. In the  $k$ -th iteration, we consider  $V^k \in \{0, 1\}^k$ . We aim to construct a map structure with (*key, prob*) pairs, where the domain of *key* is

<sup>2</sup>Note that the reason why  $A_0(V) \neq A_1(\bar{V})$  when  $u(V) = w(V)$  is that as  $0.5 \cdot e^{u(V)} = 0.5 \cdot e^{w(V)}$ , based on Theorem 1,  $BV(V) = BV(\bar{V}) = 0$ , so  $A_0(V) = 0.5 \cdot e^{u(V)}$  and  $A_1(\bar{V}) = 0$ .



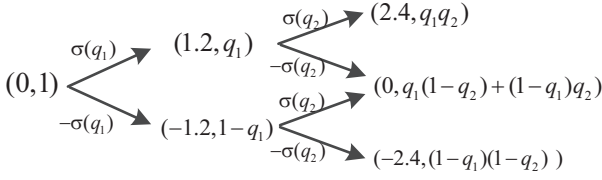


Figure 4: The iterative approach  $\{ \text{pair}(\text{key}, \text{prob}) \}$

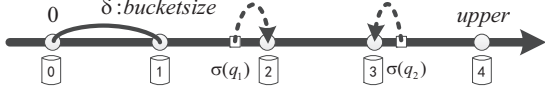


Figure 5: Principle of the bucket array.

$\{ R(V^k) \mid V^k \in \{0, 1\}^k \}$ , and the corresponding value of the *key*, or *prob* is

$$\text{prob} = \sum_{R(V^k)=\text{key} \wedge V^k \in \{0, 1\}^k} e^{u(V^k)}. \quad (7)$$

Suppose in the  $k$ -th iteration, such a map structure is constructed. Then in the next iteration, we can generate a new map structure from the old map structure: for each  $(\text{key}, \text{prob})$  in the old map structure, based on the possible choices of  $v_{k+1}$  and by considering two formulas in Equation 6, we have

1. for  $v_{k+1} = 0$ , the new key  $\text{key} + \sigma(q_{k+1})$  is generated and  $\text{prob} \cdot q_{k+1}$  is added to the prob of the new key;
2. for  $v_{k+1} = 1$ , the new key  $\text{key} - \sigma(q_{k+1})$  is generated and  $\text{prob} \cdot (1 - q_{k+1})$  is added to the prob of the new key.

EXAMPLE 4. We give an example to illustrate the above process in Figure 4, where  $n = 2$  and  $\sigma(q_1) = \sigma(q_2) = 1.2$ . Starting from  $(0, 1)$ , for  $v_1 = 0$  and  $v_1 = 1$ , it respectively creates  $(\sigma(q_1) : q_1)$  and  $(-\sigma(q_1) : (1 - q_1))$  in the first iteration. Then it leverages the stored  $(\text{key}, \text{prob})$  pair to generate new pairs in the second iteration by considering different  $v_2$ . Note that as  $\sigma(q_1) = \sigma(q_2)$ , if  $(\sigma(q_1), q_1)$  takes  $v_2 = 1$  and  $(-\sigma(q_1), (1 - q_1))$  takes  $v_2 = 0$ , then they go to the same key  $= 0$ , and their new prob,  $q_1 \cdot (1 - q_2)$  and  $(1 - q_1) \cdot q_2$  are added together.

### 4.3 Bucket-Based Approximation Algorithm

By our intractability result for JQ we know that the domain of keys, or  $\{R(V) \mid V \in \{0, 1\}^n\}$  is exponential. In order to address this issue, we set a controllable parameter  $\text{numBuckets}$  and map  $\sigma(q_i)$  to a bucket integer  $b_i \in [0, \text{numBuckets}]$ , where the interval between adjacent buckets, called *bucket size* (denoted as  $\delta$ ) is the same. Suppose  $\text{numBuckets} = d \cdot n$ , i.e., a constant multiple of the number of jury members, then, for each iteration, the number of possible values in the *key* is bounded by  $2dn^2 + 1$  (in the range  $[-dn^2, dn^2]$ ) Considering all  $n$  iterations, the time complexity is bounded by  $\mathcal{O}(dn^3)$ , which is of polynomial order.

We detail this process in Algorithm 1. To start with, the function `GetBucketArray` assigns  $b_i$  to worker  $j_i$  based on  $\sigma(q_i)$ . The computation of  $b_i$  proceeds as follows. At first, we fix a range  $[0, \text{upper}]$  where  $\text{upper} = \max_{i \in [1, n]} \{\sigma(q_i)\}$ . Then, we divide the range into  $\text{numBuckets}$  of buckets with equal length, denoted by  $\delta = \frac{\text{upper}}{\text{numBuckets}}$ . Finally, each worker  $j_i$ 's bucket number  $b_i$  is assigned to its closet bucket:  $b_i = \left\lceil \frac{\sigma(q_i)}{\delta} - \frac{1}{2} \right\rceil$ . Figure 5 illustrates an example where  $\text{numBuckets} = 4$ . Since  $\sigma(q_1)$  is the closet to bucket number 2, so  $b_1 = 2$ , and similarly  $b_2 = 3$ .

---

#### Algorithm 1 EstimateJQ

---

**Input:**  $J = \{j_1, j_2 \dots j_n\}, \text{numBuckets}, n$

**Output:**  $\widehat{JQ}$

```

1:  $b = \text{GetBucketArray}(J, \text{numBuckets}, n)$ ;
2:  $b = \text{Sort}(b)$ ; // sort in decreasing order, for pruning
3:  $J = \text{Sort}(J)$ ; // sort based on worker quality, similar as above
4:  $\text{aggregate} = \text{AggregateBucket}(b, n)$ ; // for pruning
5:  $\widehat{JQ} = 0$ ; // estimated JQ
6:  $SM[0] = 1$ ; // initialize a map structure
7: for  $i = 1$  to  $n$  do
8:    $M = \text{map}()$ ; // initialize an empty map structure
9:   for  $(\text{key}, \text{prob}) \in SM$  do
10:     $\text{flag}, \text{value} = \text{Prune}(\text{key}, \text{prob}, \text{aggregate}[i])$ ;
11:    if  $\text{flag} = \text{true}$  then
12:       $\widehat{JQ} += \text{value}$ ;
13:    continue // for pruning
14:    if  $\text{key} + b[i] \notin M$  then
15:       $M[\text{key} + b[i]] = 0$ ;
16:       $M[\text{key} + b[i]] += \text{prob} \cdot q_i$ ; // for  $v_i = 0$ 
17:    if  $\text{key} - b[i] \notin M$  then
18:       $M[\text{key} - b[i]] = 0$ ;
19:       $M[\text{key} - b[i]] += \text{prob} \cdot (1 - q_i)$ ; // for  $v_i = 1$ 
20:     $SM = M$ ;
21: for  $(\text{key}, \text{prob}) \in SM$  do
22:   if  $\text{key} > 0$  then
23:      $\widehat{JQ} += \text{prob}$ ;
24:   if  $\text{key} = 0$  then
25:      $\widehat{JQ} += 0.5 \cdot \text{prob}$ ;
26: return  $\widehat{JQ}$ ;
```

---



---

#### Algorithm 2 Pruning Techniques

---

```

def AggregateBucket( $b, n$ ):
   $\text{aggregate} = [0, 0 \dots 0]$ ; // n elements, all 0
  for  $i = n$  to  $1$  do
    if  $i = n$  then
       $\text{aggregate}[i] = b[i]$ ;
    else
       $\text{aggregate}[i] = \text{aggregate}[i + 1] + b[i]$ ;
  return  $\text{aggregate}$ 

def Prune( $\text{key}, \text{prob}, \text{number}$ ):
   $\text{flag} = \text{false}$ ;
  if  $\text{key} > 0$  and  $\text{key} - \text{number} > 0$  then
     $\text{flag} = \text{true}$ ;  $\text{value} = \text{prob}$ ;
  if  $\text{key} < 0$  and  $\text{key} + \text{number} < 0$  then
     $\text{flag} = \text{true}$ ;  $\text{value} = 0$ ;
  return  $\text{flag}, \text{value}$ ;
```

---

After mapping each worker to a bucket  $b_i$ , we iterate over  $n$  workers (step 7-20). For a given worker  $j_i$ , based on each  $(\text{key}, \text{prob})$  pair in the stored map  $SM$ , we update *key* and *prob*, based on two possible values of  $v_i$  (steps 14-19)<sup>3</sup> in the new map  $M$ .  $SM$  will then be updated as the newly derived map  $M$  for next iteration (step 20). Finally, the  $(\text{key}, \text{value})$  pairs in  $SM$  are used in the evaluation of the Jury Quality (steps 21-25), based on the cases in Figure 3.

**Pruning Techniques.** We can further improve the running time of the approximation algorithm by applying some pruning techniques in Algorithm 2, in order to prune redundant computations. For example, assume  $n = 5$ , and the derived  $b = [3, 7, 4, 3, 2]$ . In the second iteration, consider the  $\text{key} = 3 + 7 = 10$  ( $v_1 = 0$  and

<sup>3</sup>Note that as we only care about the sign (+, 0 or -) of  $R(V)$ , and we approximate  $\sigma(q_i)$  as  $\delta \cdot b_i$ , we can map  $\sigma(q_i)$  to  $b_i$  and add/subtract the integer  $b_i$ .

$v_2 = 0$ ). No matter what the rest of the three votes are, the aggregated buckets cannot be negative (since  $4 + 3 + 2 = 9 < 10$ ), so we can safely prune the search space for  $key = 10$  (which takes  $2^3 = 8$  computations). To further increase the efficiency, in Algorithm 2 we first sort the bucket array and  $J$  in decreasing order (step 2-3), guaranteeing that the highest bucket is considered first, and then compute the *aggregate* array via `AggregateBucket` (step 4), which makes the pruning phase (step 10-13) more efficient. The function `Prune` uses *aggregate* to decide whether to prune or not.

#### 4.4 Approximation Error Bound

Let  $\widehat{JQ}$  denote the estimated value returned by Algorithm 1, and  $JQ$  denote the real Jury Quality. We evaluate the *additive error* bound on  $|JQ - \widehat{JQ}|$  and we can prove that:

$$\widehat{JQ} \leq JQ \quad \text{and} \quad JQ - \widehat{JQ} < e^{\frac{n\delta}{4}} - 1, \quad (8)$$

where  $n$  is the jury size and  $\delta = \frac{upper}{d \cdot n}$  is the bucket size. Interested readers can refer to technical report [15] for the detailed proof.

We next show that the bound is very small ( $< 1\%$  by setting  $d \geq 200$ ) in real cases. First we notice that (i)  $\sigma(q)$  is a strictly increasing function and (ii)  $\sigma(0.99) < 5$ . So let us assume  $upper < 5$ . We can safely make the assumption, since if not, there exists a worker of quality  $q_i > 0.99$ , and then  $JQ \in (0.99, 1]$ , as Lemma 1 will show. Thus we can just return  $\widehat{JQ} = q_i > 0.99$ , which makes  $JQ - \widehat{JQ} < 1\%$ . After dividing the interval  $[0, upper]$  into  $d \cdot n$  equal buckets, we have  $\delta < \frac{5}{d \cdot n}$ . Using this  $\delta$  bound in Equation 8, we have  $JQ - \widehat{JQ} < e^{\frac{5}{4 \cdot d}} - 1$ . By setting  $d \geq 200$ , the bound is  $JQ - \widehat{JQ} < 0.627\% < 1\%$ .

#### 4.5 Incorporation of Prior

In the previous section, we have assumed a prior  $\alpha = 0.5$ . Here, we drop this assumption and show how we can adapt our approaches to a generalized prior  $\alpha \in [0, 1]$ , given by the task provider. By Theorem 3, it turns out this is equivalent to computing  $JQ(J', BV, 0.5)$ , where  $J'$  is obtained by adding a worker (with quality  $\alpha$ ) to  $J$ :

**THEOREM 3.** *Given  $\alpha$  and  $J$ ,  $JQ(J, BV, \alpha) = JQ(J', BV, 0.5)$ , where  $J' = J \cup \{j_{n+1}\}$  and  $q_{n+1} = \alpha$ .*

Due to lack of space, interested readers can refer to technical report [15] for the detailed proof.

Thus we can use Algorithm 1 for any prior  $\alpha$ , by adding to the jury a pseudo-worker of quality  $\alpha$ . Moreover, the approximation error bound proved in Section 4.4 also holds.

In summary, to compute  $JQ(J, BV, \alpha)$ , we have developed an approximation algorithm with time complexity  $\mathcal{O}(d \cdot n^3)$ , with an additive error bound within  $1\%$ , for  $d \geq 200$ .

### 5. JURY SELECTION PROBLEM (JSP)

Now we focus on addressing  $J^* = \arg \max_{J \in \mathcal{C}} JQ(J, BV, \alpha)$ , for  $\mathcal{C}$ , the set of all feasible juries (i.e.,  $\mathcal{C} = \{J \mid J \subseteq W \wedge \sum_{i=1}^n c_i \leq B\}$ ).

Before formally addressing JSP, we turn our attention to two monotonicity properties of  $JQ(J, BV, \alpha)$ : with respect to varying the jury size ( $|J|$ ), and with respect to a worker  $j_i$ 's quality ( $q_i$ ). These properties can help us solve JSP under certain cost constraints.

**LEMMA 1 (MONOTONICITY ON JURY SIZE).** *Given  $\alpha$  and  $J$ ,  $JQ(J, BV, \alpha) \leq JQ(J', BV, \alpha)$ , where  $J' = J \cup \{j_{n+1}\}$ .*

**LEMMA 2 (MONOTONICITY ON WORKER QUALITY).**

*Given  $\alpha$  and  $J$ . Let  $J' = J$  except that  $q'_{i_0} \geq q_{i_0} \geq 0.5$  for some  $i_0$ , then  $JQ(J', BV, \alpha) \geq JQ(J, BV, \alpha)$ .*

**PROOF.** Due to space limits, interested reader can refer to technical report [15] about the proofs for Lemma 1 and 2.  $\square$

A direct consequence of Lemma 1 is that “the more workers, the better JQ for BV”. So for the case that each worker will contribute voluntarily ( $c_i = 0$  for  $1 \leq i \leq N$ ) or the budget constraint satisfies on all subsets of the candidate workers  $W$  (i.e.,  $B \geq \sum_{i=1}^N c_i$ ), we can select all workers in  $W$ .

Lemma 2 shows that a worker with higher quality contributes not less in JQ compared with a lower quality worker. For the case that each worker has the same cost requirement  $c$ , i.e.,  $c_i = c_j = c$  for  $i, j \in [1, N]$ , we can select the top- $k$  workers sorted by their quality in decreasing order, where  $k = \min\{\lfloor \frac{B}{c} \rfloor, N\}$ .

Although the above two properties can indicate us to solve JSP under certain conditions, the case for JSP with arbitrary individual cost is much more complicated as we have to consider not only the worker  $j_i$ 's quality  $q_i$ , but also her cost  $c_i$ , and both may vary between different workers.

We can formally prove that JSP is NP-hard in Theorem 4. Note that JSP, in general, is NP-hard due to the fact that it cannot avoid computing  $JQ(J, BV, \alpha)$  at each step, which is an NP-hard problem itself. Moreover, even if we assume the existence of a polynomial oracle for computing  $JQ(J, BV, \alpha)$  (e.g., Algorithm 1), the problem still remains NP-hard, as we can reduce a  $n$ -th order Knapsack Problem [7] to it. Interested readers can refer to the technical report [15] for more details.

**THEOREM 4.** *Solving JSP is NP-hard.*

#### 5.1 Heuristic Solution

To address the computational hardness issue, we use the *simulated annealing heuristic* [19], which is a stochastic local search method for discrete optimization problems. This method can escape local optima and is proved to be effective in solving a variety of computationally hard problems [5, 10].

The simulated annealing heuristic mimics the cooling process of metals, which converge to a final, “frozen” state. A temperature parameter  $T$  is used and iteratively reduced until it is small enough. For a specific value of  $T$ , the heuristic performs several local neighbourhood searches. There is an objective value on each location, and let  $\Delta$  denote the difference in objective value between the searched location and the original location. For each local search, the heuristic makes a decision whether to “move” to the new location or not based on  $T$  and  $\Delta$ :

1. if the move will not decrease the objective value (i.e.,  $\Delta \geq 0$ ), then the move is accepted;
2. if the move will decrease the objective value (i.e.,  $\Delta < 0$ ), the move is accepted with probability  $\exp(-\frac{\Delta}{T})$ , i.e., by sampling from a Boltzmann distribution [21].

The reason for not immediately rejecting the move towards a worse location is that it tries to avoid getting stuck in local optima. Intuitively, when  $T$  is large, it is freer to move than at lower  $T$ . Moreover, a large  $\Delta$  restricts the move as it increases the chances of finding a very bad case.

We can apply the simulated annealing heuristic to solve JSP in Algorithm 3 by assuming that each location is a jury set  $J \subseteq W$  and its objective value is  $JQ(J, BV, \alpha)$ . What is important in simulated annealing is the design of local search. Before introducing

---

**Algorithm 3** JSP

---

**Input:**  $W = \{j_1, j_2, \dots, j_N\}, B, N$   
**Output:**  $\hat{J}$

- 1:  $T = 1.0$ ; // initial temperature parameter
- 2:  $X = [x_1 = 0, x_2 = 0, \dots, x_N = 0]$ ; // all initialized as 0
- 3:  $\hat{J} = \emptyset$ ; // estimated optimal jury set  $J^*$
- 4:  $M = 0$ ; // the overall monetary incentive for selected workers
- 5:  $H = \emptyset$ ; // the set containing indexes for selected workers
- 6: **while**  $T \geq \epsilon$  **do**
- 7:   **for**  $i = 1$  to  $N$  **do**
- 8:     randomly pick an index  $r \in \{1, 2, \dots, N\}$ ;
- 9:     **if**  $x_r = 0$  and  $M + c_r \leq B$  **then**
- 10:        $x_r = 1$ ;  $M = M + c_r$ ;
- 11:        $\hat{J} = \hat{J} \cup \{j_r\}$ ;  $H = H \cup \{r\}$ ;
- 12:     **else**
- 13:        $X, M, \hat{J}, H = \text{Swap}(X, M, \hat{J}, H, r, B, N)$ ;
- 14:      $T = T/2$ ; // cool the temperature
- 15: **return**  $\hat{J}$ ;

---

---

**Algorithm 4** Swap

---

**Input:**  $X, M, \hat{J}, H, r, B, N$   
**Output:**  $X, M, \hat{J}, H$

- 1: **if**  $x_r = 0$  **then**
- 2:   randomly pick an index  $k \in H$ ;
- 3:    $a = k$ ;  $b = r$ ; // store the index
- 4: **else**
- 5:   randomly pick an index  $k \in \{1, 2, \dots, N\} \setminus H$ ;
- 6:    $a = r$ ;  $b = k$ ; // store the index
- 7: **if**  $M - c_a + c_b \leq B$  **then**
- 8:    $\Delta = \text{EstimateJQ}(\hat{J} \setminus \{j_a\} \cup \{j_b\}) - \text{EstimateJQ}(\hat{J})$ ;
- 9:   **if**  $\Delta \geq 0$  **or**  $\text{random}(0, 1) \leq \exp(-\frac{\Delta}{T})$  **then**
- 10:      $x_a = 0$ ;  $x_b = 1$ ;  $M = M - c_a + c_b$ ;
- 11:      $\hat{J} = \hat{J} \setminus \{j_a\} \cup \{j_b\}$ ;  $H = H \setminus \{a\} \cup \{b\}$ ;
- 12: **return**  $X, M, \hat{J}, H$

---

our design of local search, we first explain some variables to keep in Algorithm 3:  $H$  is used to store the indexes of selected workers,  $M$  is used to store their aggregated cost, and  $X = [x_1, x_2, \dots, x_N]$  is used to keep the current state of each worker ( $x_i = 1$  indicates that worker  $j_i$  is selected and 0 otherwise). Starting from an initial  $X$ , we iteratively decrease  $T$  (step 14) until  $T$  is small enough (step 6). In each iteration, we perform  $N$  local searches (steps 7-13), by randomly picking an index  $r$  out of the  $N$  worker indexes. Based on the randomly picked  $x_r$ , we either select the worker if adding the worker does not violate the budget  $B$  (steps 9-11), or execute **Swap**, which is described in Algorithm 4. The decision to swap is made based on different  $x_r$  values:

1. if  $x_r = 0$ , a randomly picked worker  $k \in H$  is replaced with worker  $r$  if the replacement does not violate the budget constraint and the move is accepted based on  $\Delta$  and  $T$ ;
2. if  $x_r = 1$ , the algorithm performs similarly to the above case, and it replaces worker  $r$  with a randomly picked worker  $k \in \{1, 2, \dots, N\} \setminus H$  if the budget constraint still satisfies and the move is accepted as above.

While the heuristic does not have any bound on the returned jury ( $\hat{J}$ ) versus the optimal jury ( $J^*$ ), we show in the experiments (Section 6) that it is close to the optimal by way of comparing the real and estimated JQ (i.e.,  $JQ(\hat{J}, BV, \alpha)$  and  $JQ(J^*, BV, \alpha)$ ).

## 6. EXPERIMENTAL EVALUATION

In this section we present the experimental evaluation of JQ and JSP, both on synthetic data and real data. For each dataset, we first

evaluate the solution to JSP first, and then give detailed analysis on the computation of JQ. The algorithms were implemented in Python 2.7 and evaluated on a 16GB memory machine with Windows 7 64bit.

### 6.1 Synthetic Dataset

#### 6.1.1 Setup

First, we describe our default settings for the experiments. Similar to the settings in [7], we generate each worker  $j_i$ 's quality  $q_i$  and cost  $c_i$  via Gaussian distributions, i.e.,  $q_i \sim \mathcal{N}(\mu, \sigma^2)$  and  $c_i \sim \mathcal{N}(\hat{\mu}, \hat{\sigma}^2)$ . We also set parameters following [7], i.e.,  $\mu = 0.7$ ,  $\sigma^2 = 0.05$ ,  $\hat{\mu} = 0.05$  and  $\hat{\sigma}^2 = 0.2$ . By default,  $B = 0.5$ ,  $\alpha = 0.5$  and the number of candidate workers in  $W$  is  $N = 50$ . For JSP (Algorithm 3), we set  $\epsilon = 10^{-8}$ ; for JQ computation (Algorithm 1), we set  $\text{numBuckets} = 50$ . To achieve statistical significance of our results, we repeat the results 1,000 times and report the average values.

#### 6.1.2 System Comparison

We first perform the comparison of JSP with previous works, in an end-to-end system experiment. Cao et al. [7] is the only related algorithm we are aware of, which solves JSP under the MV strategy in an efficient manner. Formally, it addresses JSP as  $\arg \max_{J \in \mathcal{C}} JQ(J, MV, 0.5)$ . We denote their system as *MVJS* (Majority Voting Jury Selection System) and our system (Figure 1) as *OPTJS* (Optimal Jury Selection System). We compare the two systems by measuring the JQ on the returned jury sets.

The results are presented in Figure 6. We first evaluate the performance of the two systems by varying  $\mu \in [0.5, 1]$  in Figure 6(a), which shows that *OPTJS* always outperforms *MVJS*, and *OPTJS* is more robust with low-quality workers. For example, when  $\mu = 0.6$ , the JQ of *OPTJS* leads that of *MVJS* for 5%. By fixing  $\mu = 0.7$ , Figure 6(b)-(d) respectively vary  $B \in [0.1, 1]$ ,  $N \in [10, 100]$ ,  $\hat{\sigma} \in [0.1, 1]$  and compare the performance of *MVJS* and *OPTJS*, which all show that *OPTJS* consistently performs better than *MVJS*. In Figure 6(b), *OPTJS* on average leads around 3% compared with *MVJS* for different  $B$ ; in Figure 6(c), *OPTJS* is better than *MVJS*, especially when the number of workers is limited (say when  $n = 10$ , *OPTJS* leads *MVJS* for more than 6%); in Figure 6(d), compared with *MVJS*, *OPTJS* is more robust with the change of  $\hat{\sigma}$ .

In summary, *OPTJS* always outperforms *MVJS* and, moreover, it is more robust with (1) lower-quality workers, (2) limited number of workers and (3) different cost variances.

#### 6.1.3 Evaluating *OPTJS*

Next, we test the approximation error of Algorithm 3 by fixing  $N = 11$  and varying  $B \in [0.05, 0.5]$ . Because of its NP-hardness,  $J^*$  is obtained by enumerating all feasible juries. We record the optimal  $JQ(J^*, BV, 0.5)$  and the returned  $JQ(\hat{J}, BV, 0.5)$  in Figure 7(a). It shows that the two curves almost coincide with each other. As mentioned in Section 6.1.1, each point in the graph is averaged over repeated experiments. Thus, we also give statistics of the difference  $JQ(J^*, BV, 0.5) - JQ(\hat{J}, BV, 0.5)$  on all the 10,000 experiments considering different  $B$  ( $B$  changes in  $[0.05, 0.5]$  with step size 0.05) in Table 3, which shows that more than 90% of them have a difference less than 0.01% and the maximum error is within 3%.

Our next experiment is to test the efficiency of Algorithm 3. We set  $B = 0.5$  and vary  $N \in [100, 500]$ . The results are shown in Figure 7(b). We observe that the running time increases linearly with  $N$ , and it is less than 2.5 seconds even for high numbers of



workers ( $N = 500$ ). It is fairly acceptable in real situations as the JSP can be done offline.

**Table 3: Counts in different error ranges**

%	[0, 0.01]	(0.01, 0.1]	(0.1, 1]	(1, 3]	(3, +∞)
Counts	9301	231	408	60	0

### 6.1.4 JQ Computation

We now turn our attention to the computation of JQ, which is an essential part of *OPTJS*. We denote here by  $n$  the jury size.

We first evaluate the optimality of BV with respect to JQ. Due to the fact the computing JQ in general is NP-hard, we set  $n = 11$  and evaluate JQ for four different strategies: two deterministic ones (MV-Majority Voting, and BV-Bayesian Voting), and two randomized ones (RBV-Random Ballot Voting<sup>4</sup> and RMV-Randomized Majority Voting). We vary  $\mu \in [0.5, 1]$  and illustrate the resulting JQ in Figure 8(a). It can be seen that the JQ for BV outperforms the others. Moreover, unsurprisingly, all strategies have their worst performance for  $\mu = 0.5$  as the workers are purely random in that case. But when  $\mu = 0.5$ , BV also performs robust (with JQ 93.3%), the reason is that other strategies are sensitive to low-quality workers, while BV can wisely decide the result by leveraging the workers’ qualities. Finally, the randomized version of MV, i.e., RMV, performs not better than MV for  $\mu \geq 0.5$ , as randomized strategies may improve the error bound in the worst case [23]. The JQ under RBV always keeps at 50% since it is purely random.

To further evaluate the performance of different strategies for different jury sizes, and for a fixed  $\mu = 0.7$ , we vary  $n \in [1, 11]$  and plot the resulting qualities in Figure 8(b). The results show that as  $n$  increases, the JQ for the two randomized strategies stay the same and BV is the highest among all strategies. To be specific, when  $n = 7$ , the BV is about 10% better than MV. In summary, BV performs the best among all strategies.

Having compared the JQ between different strategies, we now focus on addressing the computation of JQ for BV, i.e.,  $JQ(J, BV, 0.5)$  in Figure 9. We first evaluate the effect of the quality variance  $\sigma^2$  with varying mean  $\mu$  in Figure 9(a). It can be seen that JQ has the highest value for a high variance when  $\mu = 0.5$ . It’s because under a higher variance, worker qualities are more likely to deviate from the mean (0.5), and so, it’s likely to have more high-quality workers.

Then we address the effectiveness of Algorithm 1 for approximating the real JQ. We first evaluate the approximation error in Figure 9(b) by varying  $numBuckets \in [10, 200]$ . As can be seen, the approximation error drops significantly with  $numBuckets$ , and is very close to 0 if we have enough buckets. In Figure 9(c) we plot the histogram of differences between the accurate JQ and the approximated JQ (or  $JQ - \widehat{JQ}$ ) over all repeated experiments by setting  $numBuckets = 50$ . It is heavily skewed towards very low errors. In fact, the maximal error is within 0.01%.

Finally, we evaluate the computational savings of the pruning techniques of Algorithm 1 by varying the number of workers  $n \in [100, 500]$  in Figure 9(d). The pruning technique is indeed effective, saving more than half the computational cost. Moreover, it scales very well with the number of workers. For example, when  $n = 500$ , the estimation of JQ runs within 2.5s without pruning technique, while finishing within 1s facilitated by the proposed pruning methods.

## 6.2 Real Dataset

<sup>4</sup>RBV randomly returns 0 or 1 with 50%.

### 6.2.1 Dataset Collection

We collected the real world data from the Amazon Mechanical Turk (AMT) platform. AMT provides APIs and allows users to batch multiple questions in Human Intelligence Tasks (HIT). Each worker is rewarded with a certain amount of money upon completing a HIT. The API also allows to set the number of assignments (denoted  $m$ ) to a HIT, guaranteeing it can be answered  $m$  times by different workers. To generate the HITs, we use the public sentiment analysis dataset<sup>5</sup>, which contains 5,152 tweets related to various companies. We randomly select 600 tweets from them, and generate a HIT for each tweet, which asks whether the sentiment of a tweet is positive or not (decision making task). The ground truth of this question is provided by the dataset. The true answers for *yes* and *no* is approximately equal, so we set the prior as  $\alpha = 0.5$ .

To perform experiments on AMT, we randomly batch 20 questions in a HIT and set  $m = 20$  for each HIT, where each HIT is rewarded \$0.02. After all HITs are finished, we collect a dataset which contains 600 decision-making tasks, and each task is answered by 20 different workers. We give several statistics on the worker answering information. There are 128 workers in total, and each of them has answered on average  $\frac{600 \times 20}{128} = 93.75$  questions. Only two workers have answered all questions and 67 workers have answered only 20 questions. We used these answers to compute every worker’s quality, which is defined as the proportion of correctly answered questions by the worker in all her answered questions. The average quality for all workers is 0.71. There are 40 workers whose qualities are greater than 0.8, and about 10% whose quality is less than 0.6.

### 6.2.2 JSP

To evaluate JSP, for each question, we form the candidate workers set  $W$  by collecting all 20 workers who answered the question, i.e., having  $N = |W| = 20$ . We follow the settings in experiments on synthetic data except that worker qualities are computed using the real-world data. We then solve JSP for each question by varying  $B \in [0.1, 1.0]$ ,  $N \in [3, 20]$  and  $\hat{\sigma} \in [0, 1]$ . We compute the average returned JQ by solving JSP for all 600 questions, which is recorded as a point in Figures 10(a)-(c), respectively. It can be seen that Figure 10(a)-(c) has a similar results pattern as Figure 6(b)-(d), i.e., experimental results on the synthetic datasets. Especially, *OPTJS* always outperforms *MVJS* in real-world scenarios.

### 6.2.3 Is JQ is a good prediction?

Finally, we try to evaluate whether JQ, defined in Definition 3, is a good way to predict the quality for BV in reality. Notice that, after workers give their votes, we can adopt BV to get the voting result, and then compare it with the true answer of the question. And thus, the goodness of BV in reality can be measured by the “accuracy”, which counts the proportion of correctly answered questions according to BV.

We now test whether JQ is a good prediction of accuracy in reality. For each question, we vary the number of votes (denoted as  $z$ ). For a given  $z \in [0, 20]$ , based on the question’s answering sequence, we collect its first  $z$  votes, then

- (i) for each question, knowing the first  $z$  workers who answered the question, we can compute the JQ by considering these workers’ qualities. Then we take the average of JQ among all 600 questions;
- (ii) by considering the first  $z$  workers’ qualities who answered the question and their votes, BV can decide the result of the question. After that, the accuracy can be computed by comparing voting result and the true answer for each question.

<sup>5</sup><http://www.sananalytics.com/lab/twitter-sentiment/>

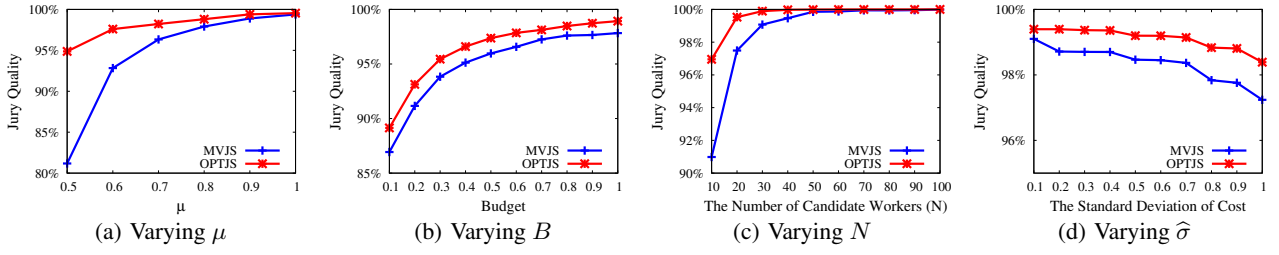


Figure 6: End-to-end system comparison

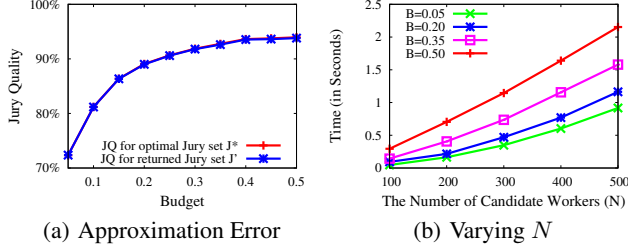


Figure 7: Efficiency & Effectiveness of *OPTJS*

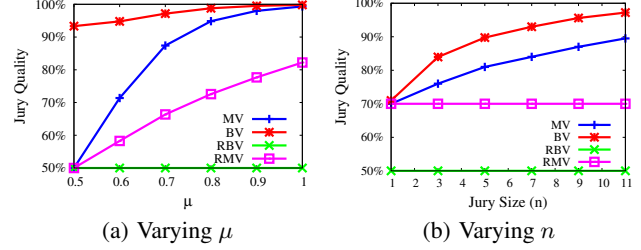


Figure 8: JQ for different strategies

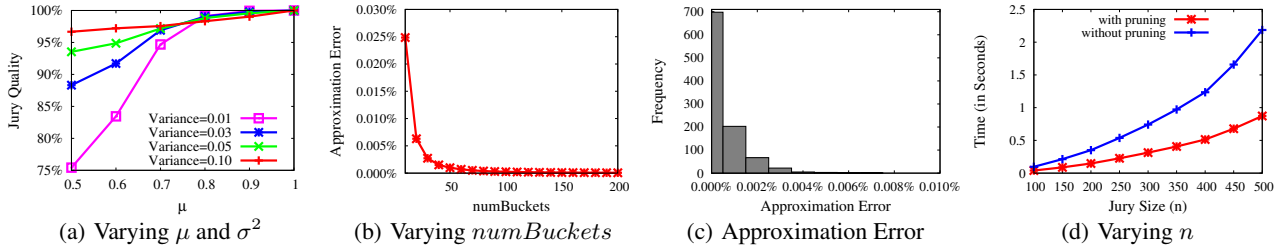


Figure 9: *JQ(J, BV, 0.5)* computation

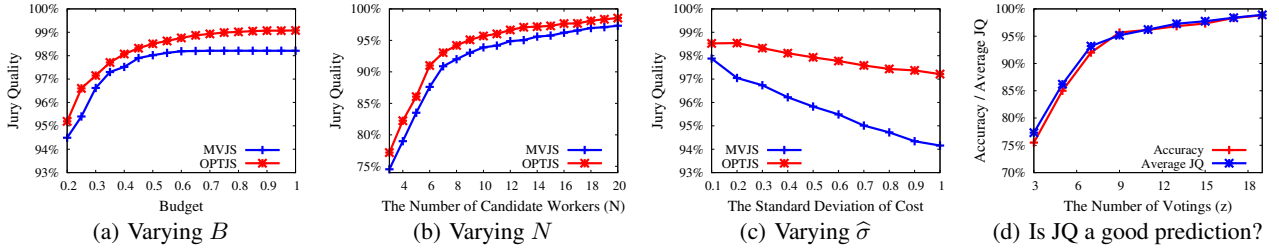


Figure 10: Real dataset evaluation

Now given a  $z \in [3, 20]$ , we compare the average JQ and accuracy in Figure 10(d), which shows that they are highly similar. Hence, it verifies that JQ for BV is really a good prediction of accuracy for BV in reality.

## 7. EXTENSIONS TO VARIOUS TASK TYPES AND WORKER MODELS

Previously we have talked about how to solve JSP under our data model. Note that we have made two assumptions: (1) it is a decision-making task with binary answer, and (2) each worker's quality is modeled as a constant. However, in reality, it is common for task provider to ask multiple-choice tasks. For example, sentiment analysis tasks [25] require workers to label the sentiment (*positive*, *neutral*, or *negative*) of each task. In addition, the worker's quality can be modeled by measuring the sensitivity and specificity of each possible answer [45], or the *confusion matrix* (CM) [18]. Specifically, a confusion matrix  $C$  is a matrix of size  $\ell \times \ell$  where each element  $C_{jk}$  encodes the probability that the worker votes for  $k$  when the true answer is  $j$ .

Our proposed algorithms can be easily extended to support these variants. Due to the space limits, we only outline our basic ideas for these extensions, and interested readers are recommended to refer to our technical report [15] for more details.

We first clarify some notations for multiple-choice task. Note that for a task with  $\ell$  possible choices, denoted as  $\{0, 1, \dots, \ell-1\}$ , and there exists one unknown true answer  $t \in \{0, 1, \dots, \ell-1\}$ . The domain of the voting from a jury  $J$  is  $\mathbf{V} \in \Omega = \{0, 1, \dots, \ell-1\}^n$ . Moreover, the prior is now a vector  $\vec{\alpha} = \{\alpha_0, \alpha_1, \dots, \alpha_{\ell-1}\}$  s.t.  $\sum_{j=0}^{\ell-1} \alpha_j = 1$ .

Following the same solution framework, we first briefly show that BV is still the optimal voting strategy with respect to JQ under this general model, and then sketch how to extend the JQ computation. Finally the extensions for JSP is addressed.

### Optimal Strategy Extension:

<sup>6</sup>For the case that each task can have multiple true answers, we can follow [30], which decomposes each task into  $\ell$  decision-making tasks, and publish these  $\ell$  tasks to workers.

To prove the optimality of BV for more general task here, we can follow the same flow as in Section 3.3. Similar to Equation 4, here  $\mathbb{E}[\mathbb{1}_{\{S(V)=t\}}]$  can be expressed as:

$$\sum_{V \in \Omega} \sum_{t=0}^{\ell-1} \Pr(\mathbf{t} = t) \cdot \Pr(V | \mathbf{t} = t) \cdot \mathbb{E}[\mathbb{1}_{\{S(V)=t\}}]. \quad (9)$$

For a given  $V \in \Omega$ , as

$$(\mathbb{E}[\mathbb{1}_{\{S(V)=0\}}], \mathbb{E}[\mathbb{1}_{\{S(V)=1\}}], \dots, \mathbb{E}[\mathbb{1}_{\{S(V)=\ell-1\}}])$$

defines a discrete probability distribution, it is not hard to prove that the optimal strategy, or  $S^*(V)$  is

$$S^*(V) = \arg \max_{t' \in \{0,1,\dots,\ell-1\}} \alpha_{t'} \cdot \Pr(V | \mathbf{t} = t'), \quad (10)$$

which corresponds to the Bayes' Theorem [3] that chooses the result as the label  $t^*$  with highest posterior probability, i.e.,  $t^* = \arg \max_{t' \in \{0,1,\dots,\ell-1\}} \Pr(\mathbf{t} = t' | V)$ . Thus  $S^* = BV$ .

### Jury Quality Computation Extension:

Recall the definition of JQ in Equation 9, to facilitate our understanding, we express  $\mathbb{E}[\mathbb{1}_{\{S(V)=t\}}]$  in the following way

$$\sum_{t'=0}^{\ell-1} \alpha_{t'} \cdot \left[ \sum_{V \in \Omega} \Pr(V | \mathbf{t} = t') \cdot \mathbb{E}[\mathbb{1}_{\{BV(V)=t'\}}] \right] \quad (11)$$

This representation enables us to consider each possible true answer separately. For each  $t' \in \{0, 1, \dots, \ell - 1\}$ , we compute

$$H(t') = \sum_{V \in \Omega} \Pr(V | \mathbf{t} = t') \cdot \mathbb{E}[\mathbb{1}_{\{BV(V)=t'\}}]$$

and then linearly combines the computed  $H(t')$  with  $\vec{\alpha}$  to get JQ. So the question falls to the computation of  $H(t')$ .

To compute  $H(t')$ , for a  $V \in \Omega$ , we have to keep track of (1) whether  $BV(V) = t'$  or not, and (2) if  $BV(V) = t'$ , the value  $\Pr(V | \mathbf{t} = t')$  should be added. Similar to the analysis in Section 4.2, we apply an iterative approach, where in each iteration, we expand  $J$  with one more worker. We develop a map structure with  $(key, prob)$  pairs to store the above two mentioned information. The *key* is an  $\ell$ -tuple

$$\left( \ln \frac{\Pr(V | \mathbf{t} = t') \cdot \alpha_{t'}}{\Pr(V | \mathbf{t} = 0) \cdot \alpha_0}, \dots, \ln \frac{\Pr(V | \mathbf{t} = t') \cdot \alpha_{t'}}{\Pr(V | \mathbf{t} = \ell - 1) \cdot \alpha_{\ell-1}} \right)$$

where the  $i$ -th element of the tuple is  $\ln \frac{\Pr(V | \mathbf{t} = t') \cdot \alpha_{t'}}{\Pr(V | \mathbf{t} = i-1) \cdot \alpha_{i-1}}$ . Intuitively, given a  $V \in \Omega$ , if  $BV(V) = t'$ , then  $\Pr(V | \mathbf{t} = t') \cdot \alpha_{t'} \geq \Pr(V | \mathbf{t} = t) \cdot \alpha_t$  for any  $t \in \{0, 1, \dots, \ell - 1\}$ , which means that the elements in the stored tuple are all  $\geq 0$ . The value *prob* corresponding to a *key* is the aggregated probability  $\Pr(V | \mathbf{t} = t')$  for the same state. In the  $k$ -th iteration, the  $V^k = \{0, 1, \dots, \ell - 1\}^k$ , for a *key*, we will generate  $\ell$  new *keys* corresponding to different votes, and update their individual *prob* for the next iteration. After  $n$  iterations, based on identifying the *keys* whose elements are all  $\geq 0$ , we can get *JQ* by aggregating the corresponding *probs*.

Since the values of elements in a tuple are unbounded, we can follow the similar idea in Section 4.3, that is to map each worker's vote to a bucket number. Note that each element in the tuple can be decomposed as the summation of individual worker's vote, thus the number of states in *keys* are bounded.

### Jury Selection Problem Extension:

To address JSP, similarly we can prove that the monotonicity on jury size by extending Lemma 1, which means that "the more workers, the better JQ" still holds for more general case. As the worker is modeled as a confusion matrix (with size  $\ell \times \ell$ ), the extension for Lemma 2 is non-trivial, and it stills remains an open question on what kind of confusion matrix will contribute more to the JQ.

Previous works [18,34] have addressed how to rank workers (or to detect spammers in all workers) based on their associated confusion matrices, which may provide good heuristics for us.

For more general cost models where each worker requires arbitrary costs, the simulated annealing heuristic regards computing JQ as a black box, so it can be simply extended here.

## 8. RELATED WORKS

**Crowdsourcing.** Nowadays, crowdsourcing has evolved as a problem solving paradigm [6] to address computer-hard tasks. To incorporate the crowd into query processing, crowdsourced databases (e.g., CrowdDB [14], Deco [31], Qurk [27] and CDAS [25]) are built, compared with traditional database systems, they do not hold the closed-world assumption. As a novel paradigm, the power of crowdsourcing has also been leveraged in other applications. For example, in Optical Character Recognition [38], Entity Resolution [39,41], Tagging [43], Schema Matching [17,44], Web Table Understanding [12], Data Cleaning [40] and so on.

**Voting Strategy.** In order to aggregate the collective wisdom of a jury, given some specific voting of a task from the jury, voting strategies are widely used to return a result, which is an estimation of the ground truth for the task. For example, Majority Voting strategy [7] strictly returns the answer corresponding to higher votes, and Random Ballot Voting [33] randomly selects the returned result. Similarly other strategies [2,9,23–25,28,29] are also talked about in a great deal. Different from their works, here we give a systematic way to classify all the strategies into two categories, and try to observe the optimal strategy in all these strategies under the Jury Selection Problem. Note that different from our problem, people may evaluate strategies under different purposes. For example, [26] analyzes the optimal Bayesian manipulation strategies by assessing the expected loss in social welfare, and [11] applies Bayesian model to take a game-theoretic approach in characterizing the symmetric equilibrium of the game with juries.

**Worker Model.** To model a worker's quality in crowdsourcing, most existing works [7,25,28,44] define it as a constant parameter indicating the probability that the worker correctly answers a question, while other work [18] defines it as a confusion matrix, which tries to capture relations between labels in questions and is specific to choices in tasks. For the methods to derive worker's quality, a normal way is to leverage the answering history. If they are not sufficient, [25] hides golden questions (questions with known ground truth) and derive the quality based on the worker's answers for them, while other work [18] applies Expectation Maximization [8] algorithm to iteratively updates worker's quality until convergence. For micro-blog services especially in Twitter, the retweet actions are usually explored to derive the error rate for each worker [7]. In our work we define worker's quality by a constant parameter (commonly used by existing works) and assume that they are known in advance. Moreover, we also extend our method to address the confusion matrix mentioned in [18].

**Online Processing.** There are also some online processing systems [4,16,25] in crowdsourcing, which addresses how to assign tasks to workers and process the workers' answers. For example, [25] proposes quality-sensitive answering model and terminate assigning questions which has got confident answers; [4] proposes an entropy-like approach to define the uncertainty of each question and assigns questions with highest uncertainty; [16] proposes cost-sensitive model to address which questions are better answered by humans or machines. Different from them, we especially evalu-

ate how to estimate the JQ before the workers are selected to answer the questions, and the quality estimation can provide statistics and guidance for the task publisher to wisely invest budget. Even though existing work [25,28] tried to estimate the quality, they assume that each worker is of the same quality.

**Expert Team Formation.** In social network, several works [13,22] studied the problem of expert team formation, that is, given the aggregated skill requirements for a task, how to find a team of experts with minimum cost (communication cost or individual financial requirement), such that the skill requirements are satisfied. Rather than the skill requirements in [13,22], we focus on the probability of drawing a correct answer, which requires to enumerate exponential number of possibilities and is indeed challenging. In fact we address the Jury Selection Problem, which is firstly proposed by [7]. But we find that the solution is sub-optimal in [7], which cannot leverage the known quality for workers. We formally address the optimal JSP problem in the paper. Some other works [9,35] also talk about how to wisely select sources for integration. The difference is that we assume the workers are given a multiple-label task and the worker model is known, while in their problem setting, the possible answers from different sources are not restricted, and the sources' exact real qualities are unknown in advance.

## 9. CONCLUSIONS

In this paper, we have studied Jury Selection Problem (JSP) for decision-making tasks, whose objective is to choose a subset of workers, such that the probability of having a correct answer (or Jury Quality, JQ) is maximized. We approach this problem from an optimality perspective. As JQ is related to voting strategy, we prove that an existing strategy, called Bayesian Voting Strategy (BV) is optimal under the JQ. Although computing JQ under BV is NP-hard, we give an efficient algorithm with theoretical guarantees. Moreover, we incorporate the task provider prior information, and we show how to extend JQ computation for different worker models and task types. Finally we evaluate JSP under BV, we prove several properties which can be used for efficient JSP computations under some constraints, and provide an approximate solution to JSP by simulated annealing heuristics.

## Acknowledgment

Reynold Cheng, Silviu Maniu, Luyi Mo, and Yudian Zheng were supported by RGC (Project HKU 711110) and HKU (Project 201311159095). We thank the reviewers for their comments.

## 10. REFERENCES

- [1] A.P.Dawid and A.M.Skene. Maximum likelihood estimation of observer error-rates using em algorithm. *Appl.Statist.*, 28(1):20–28, 1979.
- [2] Ashish Goel and David Lee. Triadic Consensus: A Randomized Algorithm for Voting in a Crowd. <http://arxiv.org/pdf/1210.0664v1.pdf>.
- [3] Bayes' Theorem. [http://en.wikipedia.org/wiki/Bayes'\\_theorem](http://en.wikipedia.org/wiki/Bayes'_theorem).
- [4] R. Boim, O. Greenspan, T. Milo, S. Novgorodov, N. Polyzotis, and W. C. Tan. Asking the right questions in crowd data sourcing. In *ICDE*, 2012.
- [5] D. Bookstaber. Simulated annealing for traveling salesman problem.
- [6] D. Braham. Crowdsourcing as a model for problem solving: An introduction and cases. *Convergence*, 14(1):75–90, 2008.
- [7] C. C. Cao, J. She, Y. Tong, and L. Chen. Whom to ask? jury selection for decision making tasks on micro-blog services. *PVLDB*, 5(11):1495–1506, 2012.
- [8] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *J.R.Statist.Soc.B*, 30(1):1–38, 1977.
- [9] X. L. Dong, B. Saha, and D. Srivastava. Less is more: Selecting sources wisely for integration. *PVLDB*, 6(2):37–48, 2012.
- [10] A. Drexler. A simulated annealing approach to the multiconstraint zero-one knapsack problem. *Computing*, 40:1–8, 1988.
- [11] J. Duggan and C. Martinelli. A bayesian model of voting in juries. *Games and Economic Behavior*, 37(2):259–294, 2001.
- [12] J. Fan, M. Lu, B. C. Ooi, W.-C. Tan, and M. Zhang. A hybrid machine-crowdsourcing system for matching web tables. In *ICDE*, 2014.
- [13] F. Farhadi, E. Hoseini, S. Hashemi, and A. Hamzeh. Teamfinder: A co-clustering based framework for finding an effective team of experts in social networks. In *ICDM Workshops*, pages 107–114, 2012.
- [14] A. Feng, M. J. Franklin, D. Kossmann, T. Kraska, S. Madden, S. Ramesh, A. Wang, and R. Xin. Crowddb: Query processing with the vldb crowd. *PVLDB*, 4(12):1387–1390, 2011.
- [15] Full version. <http://i.cs.hku.hk/~ydzhang2/Jury/Jury.pdf>.
- [16] J. Gao, X. Liu, B. C. Ooi, H. Wang, and G. Chen. An online cost sensitive decision-making method in crowdsourcing systems. In *SIGMOD*, 2013.
- [17] N. Q. V. Hung, N. T. Tam, Z. Miklós, and K. Aberer. On leveraging crowdsourcing techniques for schema matching networks. In *Database Systems for Advanced Applications*, pages 139–154. Springer, 2013.
- [18] P. Ipeirotis, F. Provost, and J. Wang. Quality management on amazon mechanical turk. In *SIGKDD workshop*, pages 64–67, 2010.
- [19] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi. Optimization by simulated annealing. *Science, New Series*, 220(4598):671–680, 1983.
- [20] A. Lacasse, F. Laviolette, M. Marchand, and F. Turgeon-Boutin. Learning with randomized majority votes. pages 162–177, 2010.
- [21] L. Landau and E. Lifshitz. *Statistical Physics. Course of Theoretical Physics 5 (3 ed.)*. Oxford: Pergamon Press, 1980.
- [22] T. Lappas, K. Liu, and E. Terzi. Finding a team of experts in social networks. In *KDD*, pages 467–476, 2009.
- [23] N. Littlestone and M. K. Warmuth. The weighted majority algorithm. *Inf. Comput.*, 108(2):212–261, Feb. 1994.
- [24] X. Liu, X. L. Dong, B. C. Ooi, and D. Srivastava. Online data fusion. *PVLDB*, 4(11):932–943, 2011.
- [25] X. Liu, M. Lu, B. C. Ooi, Y. Shen, S. Wu, and M. Zhang. Cdas: A crowdsourcing data analytics system. *PVLDB*, 5(10):1040–1051, 2012.
- [26] T. Lu, P. Tang, A. D. Procaccia, and C. Boutilier. Bayesian vote manipulation: Optimal strategies and impact on welfare. In *UAI*, pages 543–553, 2012.
- [27] A. Marcus, E. Wu, S. Madden, and R. C. Miller. Crowdsourced databases: Query processing with people. In *CIDR*, pages 211–214, 2011.
- [28] L. Mo, R. Cheng, B. Kao, X. S. Yang, C. Ren, S. Lei, D. W. Cheung, and E. Lo. Optimizing plurality for human intelligence tasks. In *CIKM*, 2013.
- [29] S. Nitzan and J. Paroush. Collective decision making and jury theorems.
- [30] A. G. Parameswaran, H. Garcia-Molina, H. Park, N. Polyzotis, A. Ramesh, and J. Widom. Crowdscreen: algorithms for filtering data with humans. In *SIGMOD*, pages 361–372, 2012.
- [31] H. Park, R. Pang, A. G. Parameswaran, H. Garcia-Molina, N. Polyzotis, and J. Widom. Deco: A system for declarative crowdsourcing. *PVLDB*, 2012.
- [32] Partition Problem. [http://en.wikipedia.org/wiki/Partition\\_problem](http://en.wikipedia.org/wiki/Partition_problem).
- [33] Random Ballot Voting. [http://en.wikipedia.org/wiki/Random\\_ballot](http://en.wikipedia.org/wiki/Random_ballot).
- [34] V. C. Raykar and S. Yu. Eliminating spammers and ranking annotators for crowdsourced labeling tasks. *Journal of Machine Learning Research*, 2012.
- [35] T. Rekatsinas, X. L. Dong, and D. Srivastava. Characterizing and selecting fresh data sources. In *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*, pages 919–930. ACM, 2014.
- [36] M. Venanzi, J. Guiver, G. Kazai, P. Kohli, and M. Shokouhi. Community-based bayesian aggregation models for crowdsourcing. In *Proceedings of the 23rd international conference on World wide web*, pages 155–164. International World Wide Web Conferences Steering Committee, 2014.
- [37] P. Venetis and H. Garcia-Molina. Quality control for comparison microtasks. In *CrowdKDD, KDD 2012 Workshop*, 2012.
- [38] L. von Ahn, B. Maurer, C. McMillen, D. Abraham, and M. Blum. Recaptcha: Human-based character recognition via web security measures. *Science*, 321(5895):1465–1468, 2008.
- [39] J. Wang, T. Kraska, M. J. Franklin, and J. Feng. Crowder: Crowdsourcing entity resolution. *PVLDB*, 5(11):1483–1494, 2012.
- [40] J. Wang, S. Krishnan, M. J. Franklin, K. Goldberg, T. Kraska, and T. Milo. A sample-and-clean framework for fast and accurate query processing on dirty data. In *International Conference on Management of Data, SIGMOD 2014, Snowbird, UT, USA, June 22-27, 2014*, pages 469–480, 2014.
- [41] J. Wang, G. Li, T. Kraska, M. J. Franklin, and J. Feng. Leveraging transitive relations for crowdsourced joins. In *Proceedings of the ACM SIGMOD International Conference on Management of Data, SIGMOD 2013, New York, NY, USA, June 22-27, 2013*, pages 229–240, 2013.
- [42] J. Whitehill, P. Ruvolo, T. Wu, J. Bergsma, and J. R. Movellan. Whose vote should count more: Optimal integration of labels from labelers of unknown expertise. In *NIPS*, pages 2035–2043, 2009.
- [43] X. S. Yang, R. Cheng, L. Mo, B. Kao, and D. W. Cheung. On incentive-based tagging. In *ICDE*, pages 685–696, 2013.
- [44] C. J. Zhang, L. Chen, H. V. Jagadish, and C. C. Cao. Reducing uncertainty of schema matching via crowdsourcing. *PVLDB*, 6(9):757–768, 2013.
- [45] B. Zhao, B. I. Rubinstein, J. Gemmell, and J. Han. A bayesian approach to discovering truth from conflicting sources for data integration. *Proceedings of the VLDB Endowment*, 5(6):550–561, 2012.