

Mean Estimation from Multiple-Choice Questions

Anson Kahng,¹ Gregory Kehne,² and Ariel D. Procaccia¹

¹Computer Science Department, Carnegie Mellon University

²Department of Mathematical Sciences, Carnegie Mellon University

Abstract

Given n values possessed by n agents, we study the problem of estimating the mean by truthfully eliciting agents' answers to multiple-choice questions about their values. We consider two natural candidates for estimation error: mean squared error (MSE) and mean absolute error (MAE). We design a randomized estimator which is asymptotically optimal for both measures in the worst case. In the case where prior distributions over the agents' values are known, we give an optimal, polynomial-time algorithm for MSE, and show that the task of computing an optimal estimate for MAE is $\#\mathcal{P}$ -hard.

1 Introduction

Organizations often desire accurate estimates of population statistics (e.g., the mean of a set of values) in settings where eliciting exact values from agents is costly or impractical. For instance, suppose that you sit on an admissions committee, and your committee's task is to accurately estimate the number of candidates who will accept their admission offer, perhaps in order to decide how many more admissions offers to extend (concretely, consider the problem of admitting students in two waves corresponding to early action and regular admission). This is a consequential problem; there are significant direct and indirect costs associated with having many more or fewer matriculants than intended.¹

Without more information either about or from the admits, this is hopeless. One potential approach is to ask each admit to provide an estimate of the probability p_i that the admit matriculates, but this is problematic because each admit may not know their own exact probability of accepting the offer, and coming up with exact probabilities places nontrivial cognitive loads on participants. Therefore, it is more reasonable to ask a multiple-choice question of the form, "How likely are you to accept this offer?" with choices "High," "Medium," and "Low." The task for the university, then, is to reconstruct an accurate estimate of the number of students who will accept their offers based on the coarse-grained information yielded by these multiple-choice queries. Specifically, the question is

"What collection of multiple-choice questions should you ask, and how should you interpret their answers so as to estimate the expected number of matriculants as accurately as possible?"

1.1 Our Approach and Results

More concretely, consider n agents, where each agent i is associated with a value $p_i \in [0, 1]$. We are interested in estimating $\|p\|_1 = \sum_{i=1}^n p_i$.

We interpret multiple-choice questions as forming a partition of $[0, 1]$ into subsets X_1, \dots, X_k , and asking agent i for the subset X_j such that $p_i \in X_j$. It is known (Lambert and Shoham 2009) that in order for multiple-choice questions to elicit truthful responses, each X_j must itself be an interval (see Section 1.2 for more details). Moreover, intervals are easier to interpret than arbitrary subsets; for example, "low" probability can be defined as $p_i \in [0, 1/3]$. Therefore, we restrict our multiple-choice questions to this framework.

Our goal is to design an estimator that consists of a set of (possibly different) multiple-choice questions which are posed to the agents, together with a function that outputs an estimate of $\|p\|_1$ based on the agents' answers to the multiple-choice questions; we denote the output of the estimator by $q(p)$. We measure the accuracy of the estimator using the *mean squared error (MSE)* $\mathbb{E}[(\|p\|_1 - q(p))^2]$ or the *mean absolute error (MAE)* $\mathbb{E}[|\|p\|_1 - q(p)|]$.²

We consider two settings corresponding to different levels of information about the agent's values. When no information about p is known (worst case), we consider the problem of designing a randomized estimator with good worst-case performance (when averaged over the estimator's randomness). We give a single randomized estimator \bar{q} which guarantees

$$\text{mse}(\bar{q}) = O\left(\frac{n}{k^2}\right), \quad \text{mae}(\bar{q}) = O\left(\frac{\sqrt{n}}{k}\right)$$

and demonstrate that this is asymptotically optimal for both measures of error.

²Throughout, we evaluate additive error with respect to estimating the sum of the p_i , from which the additive error with respect to mean estimation can easily be derived.

¹In fact, our work on this paper originated from thinking about ways to address this problem in our institution.

In the second setting, each p_i is drawn from a known distribution P_i ; we consider the problem of designing a deterministic estimator which performs well on average (over the randomness of the p_i). We present an MSE-optimal estimator, and show that the problem of devising an MAE-optimal estimator is $\#P$ -hard.

1.2 Related Work

Mechanism design for information elicitation has been widely studied in computer science and economics (Zohar and Rosenschein 2008; Chen and Kash 2011; Waggoner and Chen 2014). One closely related paper by Lambert and Shoham (2009) studies the problem of eliciting truthful answers to multiple-choice questions. In their framework, payments based on agent reports and the observed outcome are necessary to induce truthful answers from agents, and the authors establish necessary and sufficient conditions on the structure of multiple-choice questions that ensure the existence of such payments. In our setting, their characterization implies that we can elicit truthful answers if and only if the questions we ask are of the form, “To which of these intervals does your p_i belong?”

More broadly, many prior works in mechanism design focus on eliciting truthful signals from agents, often through direct verification mechanisms like strictly proper scoring rules (Gneiting and Raftery 2007; Brier 1950; Good 1952; Winkler 1969) and prediction markets (Wolfers and Zitzewitz 2004; Berg et al. 2008). In a related vein, Radanovic and Faltings (2014) developed a mechanism for truthful elicitation of continuous signals, but we consider the problem of reconstructing a continuous value from discrete reports.

Additionally, Soloviev and Halpern (2018) consider the problem of acquiring information with resource limitations, where budget constraints on the number of tests in their setting roughly map to constraints in our setting on the granularity of queries. However, their setting involves noisy tests of Boolean formula truth values, as opposed to estimating a population statistic.

Furthermore, by considering estimation error as it varies over a range of k , we can investigate the relationship between the elicitation of values and the accuracy of our estimate, a tradeoff which has been studied for intelligent decision-making systems in a sequential setting by Boutilier (2002).

Alternatively, the task of optimally estimating $\sum_i p_i$ in terms of MSE or MAE can be viewed as variants of k -means and k -medians clustering, respectively. On the one hand, it resembles a special case of clustering in that P is product distribution over $[0, 1]^n$ (as opposed to a general distribution over a metric space), and our ‘clusters’ C correspond to vectors p which yield the same vector of multiple-choice answers, and the C are constrained to have a product structure as well. On the other hand, it is distinct from clustering in that r reports scalar representative ℓ_1 norms for the clusters, and so we take the ℓ_1 norm of our C before calculating error.

This line of inquiry is also related to the notion of approximate query processing (AQP) in the data mining and databases literature, which is the practice of answering expensive aggregation queries with limited resources.

Multiple research groups have focused on the bounded-error estimation of aggregates (e.g., sums of values in a database), but mostly through sampling techniques as opposed to summarization approaches (Jagadish et al. 1998; Chaudhuri, Das, and Narasayya 2007; Chaudhuri et al. 2001; Babcock, Chaudhuri, and Das 2003). More recently, there has been some work on summarizing data distributions with histograms in order to minimize the ℓ_2 distance between the distribution and the histogram approximation (Acharya et al. 2015), but this only coincides with our MSE setting when we query exactly one agent.

2 The Model

We consider a set of agents $[n] = \{1, \dots, n\}$, each with an associated number $p_i \in [0, 1]$.

Our goal is to devise a scheme for estimating the sum $\|p\|_1 = \sum_{i=1}^n p_i$ (or, equivalently, the mean) to minimize additive error. We may ask each agent i which of k intervals contains their p_i , and so our estimator chooses n partitions $B_i := \{B_{i,1}, \dots, B_{i,k}\}$ of $[0, 1]$ into intervals, and for each i the function $b_i : [0, 1] \rightarrow [k]$ poses the question to agent i and returns their response; $b_i(p_i) = j$ if $p_i \in B_{i,j}$. We refer to $b(p) := (b_1(p_1), \dots, b_n(p_n))$ as the *classifier*. Next the *aggregator* $r : [k]^n \rightarrow \mathbb{R}$ takes the agents’ responses and estimates $\sum_i p_i$; we refer to the output of r as the *report*.

Let $c_i := b_i(p_i)$ denote the answer of each agent. For $c \in [k]^n$, the box $C_c := \prod_i B_{i,c_i} = b^{-1}(c)$ is the set of (p_1, \dots, p_n) for which each agent i answers c_i . In terms of the classifier and aggregator, our goal may be restated as finding the *estimator*

$$q := r \circ b : [0, 1]^n \rightarrow \mathbb{R}$$

that minimizes expected error.

We consider two natural measures of error, mean squared error (MSE) and mean absolute error (MAE), which are defined to be

$$\begin{aligned} \text{mse}(q) &= \mathbb{E} [(\|p\|_1 - q(p))^2], \\ \text{mae}(q) &= \mathbb{E} [\|p\|_1 - q(p)]. \end{aligned}$$

When the P_i are adversarially chosen, these expectations are taken over the randomness of the estimator, and when each p_i is drawn from a known distribution P_i , these expectations are taken over the product distribution P .

Finally, throughout the paper we denote the centroid of $C \subset \mathbb{R}^n$ by $\mu(C)$. More formally,

$$\mu(C) := \frac{1}{P(C)} \int_C p \, dP,$$

where P is a measure on \mathbb{R}^n .

3 Worst-Case Guarantees

In this section we consider the case where no knowledge of the p_i is assumed, and establish upper and lower bounds on the performance of deterministic and randomized estimators.

First, suppose that the estimator $q = r \circ b$ is deterministic. For fixed b , it is clear that r should report the sum corresponding to the center point of each box $C_c = \prod_i B_{i,c_i}$,

since this minimizes the worst-case error across all $p \in C_c$. Accordingly, an adversary will seek the box with the largest ℓ_1 diameter, which can be identified by finding the $B_{i,j}$ of maximum diameter for each i . Therefore the worst-case optimal deterministic estimator chooses equipartitions $\mathcal{B}_i = \{[0, \frac{1}{k}], \dots, [\frac{k-1}{k}, 1]\}$, reports the ℓ_1 norm of the center of each box, and satisfies

$$\begin{aligned} \max_p (\|p\|_1 - q(p))^2 &= \frac{n^2}{4k^2} \\ \max_p \|\|p\|_1 - q(p)\| &= \frac{n}{2k}, \end{aligned}$$

and this is clearly tight. This is the uniform estimator, and we will denote it $q_U := r_U \circ b_U$, where b_U partitions each $[0, 1]$ into equal-size subintervals, and $r_U(c)$ is the ℓ_1 norm of the center of each C_c .

A randomized estimator, however, can perform significantly better over worst-case inputs. Indeed, consider the following randomized estimator, which we denote by $\bar{q} = r \circ \bar{b}$. We construct a randomized classifier by choosing “shifts” $s_i \in [0, \frac{1}{k-1})$ for each i uniformly and independently. Take

$$\bar{b}_i(p_i) := j \quad \text{s.t.} \quad \frac{j-1}{k-1} \leq p_i + s_i < \frac{j}{k-1}.$$

Intuitively, this partitions $[0, 1]$ into k subintervals by taking the $k-1$ thresholds $1/(k-1), \dots, 1$ and shifting them left by s_i . Then make the (deterministic) reports

$$r_i(j) := \frac{j-1}{k-1}$$

for $j \in [k]$, and define $\bar{b}(q) := (\bar{b}_1(p_1), \dots, \bar{b}_n(p_n))$ and take the aggregator $r(c) := r_1(c_1) + \dots + r_n(c_n)$. Putting these together yields the randomized estimator

$$\bar{q} := r \circ \bar{b},$$

which is illustrated in Figure 1.

Our main result for this section is the following theorem.

Theorem 1. *In the worst-case setting, the randomized estimator \bar{q} satisfies*

$$\begin{aligned} \text{mse}(\bar{q}) &= O(n/k^2) \\ \text{mae}(\bar{q}) &= O(\sqrt{n}/k). \end{aligned}$$

Moreover, these bounds are asymptotically optimal for both measures of error.

The rest of the section is devoted to proving the theorem. We do so via several lemmas, starting with the upper bound.

Lemma 1. *The randomized estimator \bar{q} satisfies*

$$\begin{aligned} \text{mse}(\bar{q}) &= O(n/k^2) \\ \text{mae}(\bar{q}) &= O(\sqrt{n}/k). \end{aligned}$$

Proof. For fixed p , we begin by analyzing the estimator coordinate-by-coordinate. Take

$$X_i := p_i - (r_i \circ \bar{b}_i)(p_i)$$

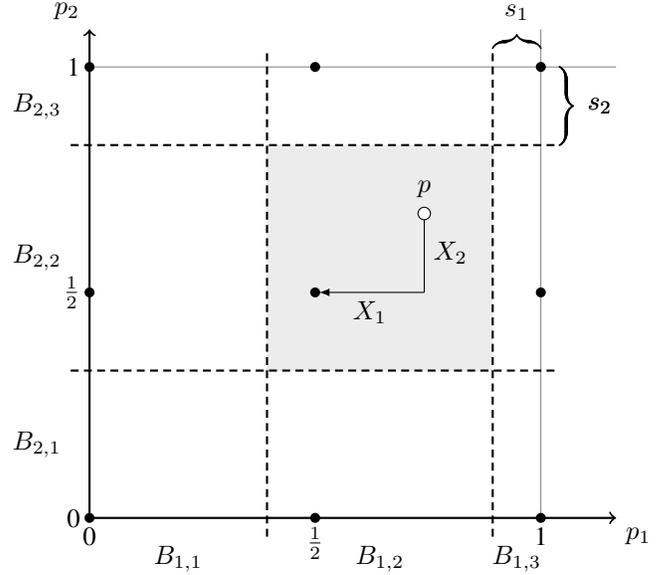


Figure 1: An illustration of \bar{q} for $n = 2$ and $k = 3$. Here $(r_1 \circ \bar{b}_1)(p_1) = 1/2$ and $(r_2 \circ \bar{b}_2)(p_2) = 1/2$, and so $\bar{q}(p) = 1/2 + 1/2 = 1$. The box $C_{(2,2)}$ is highlighted.

to be the (signed) error of \bar{q} in coordinate i , and let $j_i := \max\{j \in [k] : \frac{j-1}{k-1} \leq p_i\}$. If

$$z_i := p_i(k-1) - (j_i - 1)$$

is the proportion of the way between multiples of $1/(k-1)$ that p_i falls, then $X_i = (z_i - 1)/(k-1)$ with probability z_i and $X_i = z_i/(k-1)$ with probability $1 - z_i$. Note that $\mathbb{E}[X_i] = 0$. By the definition of r and the independence of the X_i ,

$$\begin{aligned} \text{mse}(\bar{q}(p)) &= \mathbb{E} \left[\left(\sum_i X_i \right)^2 \right] = \sum_i \text{Var} [X_i] \\ &= \sum_i \frac{z_i^2(1-z_i) + z_i(1-z_i)^2}{(k-1)^2} \leq \frac{n/4}{(k-1)^2} \\ &= O\left(\frac{n}{k^2}\right). \end{aligned}$$

We now turn to the MAE case. Note that X_i is bounded in some range of width $1/(k-1)$, and that

$$\text{mae}(\bar{q}(p)) = \mathbb{E}[\|p\|_1 - \bar{q}(p)] = \mathbb{E} \left[\left| \sum_i X_i \right| \right].$$

Next we apply Hoeffding’s inequality in order to upper bound this expectation. By Hoeffding (1963),

$$\Pr \left[\left| \sum_i X_i \right| \geq t \right] \leq 2 \exp \left\{ \frac{-2t^2}{n} \right\},$$

and so choosing $t = m\sqrt{n}/(k-1)$ for $m \in \mathbb{N}$ yields

$$\Pr \left[\left| \sum_i X_i \right| \geq \frac{m\sqrt{n}}{(k-1)} \right] \leq \frac{2}{e^{2m^2}}. \quad (1)$$

Finally, let $X = \sum_i X_i$ and observe that for any $\sigma > 0$,

$$\begin{aligned} \text{mae}(\bar{q}(p)) &= \mathbb{E}[|X|] \\ &\leq \sum_{m=1}^{\infty} \sigma m \Pr[|X| \in [\sigma(m-1), \sigma m]] \\ &\leq \sum_{m=1}^{\infty} \sigma m \Pr[|X| \geq \sigma(m-1)]. \end{aligned}$$

Taking $\sigma = \sqrt{n}/(k-1)$ and applying Equation (1),

$$\leq \frac{\sqrt{n}}{k-1} \sum_{m=1}^{\infty} \frac{2m}{e^{2(m-1)^2}}.$$

This infinite series converges, which implies that $\text{mae}(\bar{q}(p)) = O(\sqrt{n}/k)$, as desired. \square

By Yao's minimax principle (Yao 1977), in order to derive a lower bound for all randomized algorithms, it suffices to fix a distribution over inputs and lower bound the average performance of any deterministic algorithm over this randomized input. To this end, we will consider the uniform distribution over $[0, 1]^n$, which we denote D , and lower bound the performance of any deterministic estimator over it. In doing so, we will prove the intuitive fact that the uniform estimator $q_U = r_U \circ b_U$ is optimal for D .

First we present a structural lemma about the optimal aggregator r for any fixed classifier b . Let $S(C, P)$ denote the probability distribution over \mathbb{R} derived by taking the ℓ_1 norm of C :

$$\Pr[S \leq x] = \Pr[\|p\|_1 \leq x | p \in C].$$

We will repeatedly make use of the following insight, which follows from calculus and reflects one of Lloyd's optimality conditions for k -means clustering (Lloyd 1982):

Lemma 2. *For $p \sim P$ and for fixed b , the MSE- and MAE-optimal aggregators $r_b^1(c)$ and $r_b^2(c)$ (respectively) report the means and medians (respectively) of $S(C_c, P)$ for all $c \in [k]^n$.*

With this in hand, we are ready to analyze the performance of q_U over D , which (by Yao's minimax principle) establishes the lower bound of Theorem 1.

Lemma 3. *If $p \sim D$ is drawn uniformly at random then the uniform estimator $q_U = r_U \circ b_U$ is optimal in terms of both MSE and MAE, and*

$$\begin{aligned} \text{mae}(q_U) &= \Omega(n/k^2), \\ \text{mse}(q_U) &= \Omega(\sqrt{n}/k). \end{aligned}$$

Proof. This follows in two steps. We will first prove that the uniform strategy $q_U = r_U \circ b_U$ is optimal for MSE and compute $\text{mse}(q_U)$ directly. Then we will prove that q_U is optimal for MAE, and finally lower bound $\text{mae}(q_U)$.

To begin, note that under the uniform distribution $P = D$ and for fixed b , Lemma 2 implies that for both MSE and MAE, the optimal aggregator is the r_b which reports the ℓ_1 norms of the centers of the C_c . Therefore we may assume that all estimators use reports r_b which are optimal for their b , and argue that b_U is the best partitioning.

The MSE case boils down to a question of variance, and it turns out we can compute the error directly. As before, let

$$X_i := \sum_{j=1}^k \mathbb{1}_{\{p_i \in B_{i,j}\}} (p_i - \mu(B_{i,j}))$$

be the (signed) error in coordinate i , the difference between their actual p_i and the center $\mu(B_{i,j})$ of the interval $B_{i,j}$ containing their p_i . Then because the p_i are independent and

$$r_b(c) = \mu(S(C_c, D)) = \sum_i \mu_D(B_{i,c_i}), \quad (2)$$

we have that for $q = r_b \circ b$:

$$\begin{aligned} \text{mse}(q) &= \mathbb{E}_D \left[(\|p\|_1 - q(p))^2 \right] = \mathbb{E}_D \left[\left(\sum_i X_i \right)^2 \right] \\ &= \text{Var} \left[\sum_i X_i \right] = \sum_i \text{Var}[X_i] = \sum_i \int_0^1 X_i^2 dx \\ &= \sum_i \sum_j \int_{B_{i,j}} (x - \mu_D(B_{i,j}))^2 dx \\ &= \sum_i \sum_j \frac{1}{12} \text{diam}(B_{i,j})^3. \end{aligned}$$

At this point, the method of Lagrange multipliers confirms that MSE is minimized when all $B_{i,j}$ are of equal diameter $1/k$, which yields precisely b_U . Therefore $q_U = r_U \circ b_U$ is optimal for D , and it has cost $\text{mse}(q_U) = n/12k^2$.

We now turn to MAE, and use a differential argument to prove that q_U is MAE-optimal when $P = D$. We begin by showing that the MAE contributed by a box C_c is convex in each of the dimensions of C_c . This will let us argue that for any classifier b with partitions $\{\mathcal{B}_i\}$ in which one is unbalanced, meaning that for some i (say $i = 1$) and some j it is the case that $\text{diam}(B_{1,j}) < \text{diam}(B_{1,j+1})$, the b^* which equalizes their widths decreases the MAE: $\text{mae}(r_{b^*} \circ b^*) < \text{mae}(r_b \circ b)$. This then implies that $q_U = r_U \circ b_U$ is MAE-optimal, since it is the only b which cannot be equalized in this way.

To see that MAE contribution is convex in each dimension of C , let $C' := \prod_{i=2}^n [-w_i, w_i]$ and consider the box $C(t) := [-t, t] \times C'$. Then by Lemma 2 the optimal report for both C and C' is 0. Let the contribution to MAE by C with report r be denoted $e(C, r)$, and call the contribution with optimal report $e(C)$. Then by symmetry,

$$\begin{aligned} e(C(t), 0) &= \int_{C(t)} \left| \sum_i p_i \right| dp \\ &= 2 \int_0^t \int_{C'} \left| x + \sum_{i=2}^n p_i \right| dp dx \\ &= 2 \int_0^t e(C', x) dx, \end{aligned}$$

and therefore

$$\frac{de(C(t), 0)}{dt} = 2e(C', t).$$

The (omitted) proof of Lemma 2 shows that $\frac{de(C,r)}{dr} > 0$ for r greater than the optimal report; we conclude that

$$\frac{d^2e(C(t), 0)}{dt^2} > 0,$$

as desired.

In order to show that b^* improves upon b , note that their boxes C_c differ only for those of the form

$$\hat{C}_c := B_{1,j} \times \prod_{i=2}^n B_{i,c_i}, \quad \tilde{C}_c := B_{1,j+1} \times \prod_{i=2}^n B_{i,c_i}$$

Pairing these up by c , it suffices to show that for each c ,

$$e(\hat{C}_c^*) + e(\tilde{C}_c^*) < e(\hat{C}_c) + e(\tilde{C}_c),$$

where C_c^* are the boxes given by b^* . This follows from the convexity of $e(C)$ in each dimension of C , established above. We conclude that q_U is MAE-optimal.

We next lower bound the MAE of q_U . Let $L \subseteq [n]$ and $H \subseteq [n]$ be the set of indices i with errors X_i that are negative and positive, respectively. It holds that

$$\begin{aligned} \text{mae}(q_U) &= \mathbb{E} \left[\left| \sum_{i \in [n]} (p_i - q_U(p_i)) \right| \right] \\ &= \mathbb{E} \left[\left| \sum_{i \in L} X_i + \sum_{i \in H} X_i \right| \right]. \end{aligned}$$

We establish that L and H are, with constant probability, of sufficiently different sizes to lead to \sqrt{n}/k error. First, note that because we are playing against a uniform adversary, the probability that each p_i is in L is $1/2$; the probability that each p_i is in H is symmetrically $1/2$. Because these are just Bernoulli random variables, applying the De Moivre-Laplace theorem (a version of the Central Limit Theorem) tells us that, as n becomes large, the sum of these Bernoulli random variables converges to a normal distribution with mean $n/2$ and standard deviation $\sqrt{n}/2$. Therefore, we know that with constant probability, the total number of agents in H is at least \sqrt{n} from its mean of $n/2$, that is,

$$\Pr[||H| - \mathbb{E}[|H|]| \geq \sqrt{n}] = \beta$$

for a constant β . It follows that, with constant probability β , $||L| - |H|| \geq 2\sqrt{n}$; denote this event by \mathcal{E} . Therefore,

$$\begin{aligned} \text{mae}(q_U) &= \mathbb{E} \left[\left| \sum_{i=1}^n X_i \right| \right] \\ &\geq \mathbb{E} \left[\left| \sum_{i=1}^n X_i \right| \middle| \mathcal{E} \right] \cdot \Pr[\mathcal{E}] \\ &= \mathbb{E} \left[\left| \sum_{i=1}^n X_i \right| \middle| \mathcal{E} \right] \cdot \beta. \end{aligned} \quad (3)$$

Now assume that \mathcal{E} occurred. Without loss of generality, assume that $|L| \leq |H|$ and, in particular, randomly break

H up into two sets, H_1 and H_2 , such that $|H_1| = |L|$ and $|H_2| \geq 2\sqrt{n}$. By construction of H_1 , the sum of the errors in indices $i \in L \sqcup H_1$ is symmetric with mean 0. It holds that

$$\begin{aligned} &\mathbb{E} \left[\left| \sum_{i=1}^n X_i \right| \middle| \mathcal{E} \right] \\ &= \mathbb{E} \left[\left| \sum_{i \in L \cup H_1} X_i + \sum_{i \in H_2} X_i \right| \middle| \mathcal{E} \right] \\ &\geq \Pr \left[\left| \sum_{i \in L \cup H_1} X_i \right| \geq 0 \middle| \mathcal{E} \right] \mathbb{E} \left[\left| \sum_{i \in H_2} X_i \right| \middle| \mathcal{E} \right] \\ &\geq \frac{1}{2} \cdot 2\sqrt{n} \cdot \frac{1}{4k} = \Omega \left(\frac{\sqrt{n}}{k} \right). \end{aligned}$$

The desired bound follows by combining this with Equation (3). \square

4 Estimation with Priors

In practice, it is useful to go beyond worst-case guarantees and ask to what degree knowing some additional information about the p_i can improve our ability to estimate their sum. Specifically, suppose that we have access to distributions P_1, \dots, P_n from which the p_i are drawn; we make the standard assumption (for computational complexity results) that these distributions are discrete. Can we design a more accurate estimator that takes prior knowledge into account?

For instance, in our running example of admissions (say for Ph.D.), one could easily come up with priors by applying machine learning to historical admissions data. The prior for a candidate would, of course, depend on relevant features such as alma mater and research interests.

It is important to note that, in this setting, the optimal estimator is deterministic, as the error of a randomized estimator is just a convex combination of deterministic estimators in its support.

4.1 An Efficient Estimator for MSE

When estimation is evaluated according to MSE, it turns out that we can answer the foregoing question in the positive:

Theorem 2. *Given discrete priors P_i for each p_i , $i \in [n]$, there is a polynomial-time estimator that is optimal with respect to MSE.*

A key component of our analysis is the following structural insight:

Lemma 4. *Given an estimator $q = r_b \circ b$, where r_b reports optimally for a given classifier b , $\text{mse}_P(q) = \sum_i \text{mse}_{P_i}(q_i)$, where $q_i = r_{b_i} \circ b_i$.*

We will prove this by employing a result from the vector quantization literature. In quantization, roughly speaking, the task is to compress some signal using only a representative subset of its values in such a way that compression error is minimized. Vector quantization performs this task for vector-valued signals, and does so using an n -dimensional partition (Gray and Neuhoff 1998). Since it typically seeks an MSE-error-minimizing vector representative

$\vec{r}(c)$ for each box C_c in its n -dimensional partition, vector quantization may be seen as an instance of k^n -means clustering in \mathbb{R}^n subject to the constraint that all clusters obey this product-of-partitions structure. In this setting, it is known (Jégou, Douze, and Schmid 2011) that if the partitions are made along independent axes, then the MSE of the optimal vector quantizer is additive:

Lemma 5. *If \vec{q} is a vector quantizer as described above which reports the centroids $\mu(C_c)$ for every C_c , and is given by $\vec{r}(c) = (r_1(c_1) \dots r_n(c_n))$ and $b(p) = (b_1(p_1), \dots, b_n(p_n))$, then*

$$\mathbb{E}_P \left[\|p - \vec{q}(p)\|_2^2 \right] = \sum_i \mathbb{E}_{P_i} \left[(p_i - r_i(b_i(p_i)))^2 \right], \quad (4)$$

where each r_i reports the centroid $\mu(B_{i,j})$.

This is a direct consequence of independence of the p_i , together with the fact that the X_i have mean 0.

The key difference between our problem and the vector quantization setting is that we measure error with respect to the aggregate $\sum_i p_i$, which means errors with respect to individual agents can “cancel out.” Nevertheless, the same structural insight holds, as we show now.

Proof of Lemma 4. We proceed in two steps. First we argue that for fixed classifier b , the optimal report r_b reports the ℓ_1 norm $r_b(c) = \|\mu_P(C_c)\|_1$ of the centroid of each box C_c . Next, we argue that this coincides with the error given on the right-hand side of Equation (4); the theorem then follows.

We begin by showing that the optimal aggregator reports $r_b(c) = \|\mu_P(C_c)\|_1$. As in Lemma 3, let $e(C, r)$ denote the contribution of $C = \prod_i B_i$ to MSE under report r . We proceed via a differential argument:

$$\begin{aligned} e(C, r) &= \int_C (\|p\|_1 - r)^2 dP \\ &= r^2 \int_C dP - 2r \int_C \|p\|_1 dP + \int_C \|p\|_1^2 dP \\ \frac{de(C, r)}{dr} &= 2r \int_C dP - 2 \int_C \|p\|_1 dP. \end{aligned}$$

Setting $\frac{de(C, r)}{dr} = 0$ yields optimal report,

$$\begin{aligned} r^* &= \frac{1}{P(C)} \int_C \|p\|_1 dP = \frac{1}{P(C)} \left(\sum_i \int_C p_i dP \right) \\ &= \frac{1}{\prod_i P_i(B_i)} \sum_i \left(\int_{B_i} p_i dP_i \prod_{j \neq i} \int_{B_j} dP_j \right) \\ &= \sum_i \frac{1}{P_i(B_i)} \int_{B_i} p_i dP_i = \sum_i \mu_{P_i}(B_i) \\ &= \|\mu_P(C)\|_1, \end{aligned}$$

where this last step is a standard property of the centroid.

Therefore for $q = r_b \circ b$, the error $\text{mse}(q)$ takes the form

$$\text{mse}_P(q) = \sum_{c \in [k]^n} \int_{C_c} (\|p\|_1 - \|\mu(C_c)\|_1)^2 dP$$

$$= \sum_{c \in [k]^n} \int_{C_c} \left(\sum_i (p_i - \mu_{P_i}(B_{i,c_i})) \right)^2 dP.$$

Since the centroid is an unbiased estimator,

$$\begin{aligned} &= \sum_{c \in [k]^n} \int_{C_c} \sum_i (p_i - \mu_{P_i}(B_{i,c_i}))^2 dP \\ &= \mathbb{E}_P \left[\|p - \vec{q}(p)\|_2^2 \right], \end{aligned}$$

as in the right-hand side of Equation (4). Therefore by Lemma 5,

$$\begin{aligned} &= \sum_i \mathbb{E}_{P_i} \left[(p_i - \mu(B_{i,c_i}))^2 \right] \\ &= \sum_i \text{mse}_{P_i}(q_i). \end{aligned}$$

□

Proof of Theorem 2. Our goal is to minimize $\text{mse}_P(q)$, and Lemma 4 implies that this can be accomplished by individually minimizing $\text{mse}_{P_i}(r_{b_i} \circ b_i)$ for each $i \in [n]$. Since the P_i are discrete distributions, finding the b_i which minimizes $\text{mse}_{P_i}(r_{b_i} \circ b_i)$ is precisely an instance of one-dimensional Euclidean k -means. It is well-known that this can be solved efficiently via dynamic programming using a recurrence described by Jensen (1969). Therefore, given priors P_i we can derive an estimator q which minimizes $\text{mse}(q)$: we first find b_1^*, \dots, b_n^* which minimize $\text{mse}_{P_i}(r_{b_i}^2 \circ b_i)$ for each $i \in [n]$, where $r_{b_i}^2$ is the optimal aggregator for b_i , then take the optimal partitions

$$b(p) = (b_1^*(p_1), \dots, b_n^*(p_n)).$$

Lemma 4 then guarantees that this $q = r_b^2 \circ b$ is optimal. □

4.2 Hardness for MAE

In contrast to the case of MSE with priors, it turns out that devising an MAE-optimal mean estimation strategy is $\#\mathcal{P}$ -hard. We show this through a pair of reductions; beginning with the counting version of KNAPSACK, which is $\#\mathcal{P}$ -complete, we show that the problem of finding a median of the sum of the (independent) Bernoulli random variables $\alpha_i \text{Bernoulli}(p)$ is $\#\mathcal{P}$ -complete. We then describe a way to derive distributions P_1, \dots, P_n from a collection of weighted Bernoulli random variables such that devising an MAE-optimal strategy from the P_1, \dots, P_n finds a median of the weighted Bernoulli sum.

We start by stating the main result of this section.

Theorem 3. *Given discrete priors P_i for each p_i , $i \in [n]$, the problem of computing an optimal estimator with respect to MAE is $\#\mathcal{P}$ -hard.*

To be concrete, the problem is defined as follows: Given a collection of discrete prior distributions P_1, \dots, P_n and a positive integer k , we are asked for a collection of partitions $b_i : [0, 1] \rightarrow [k]$ and an aggregator $r : [k]^n \rightarrow \mathbb{R}$ which together minimize

$$\text{mae}(r \circ b) = \mathbb{E}_P [\| \|p\|_1 - r \circ b(p) \|],$$

where $P = \prod_i P_i$.

To prove the theorem, we will reduce from the following problem. Given a rational $x \in [0, 1]$ and nonnegative integer weights $\alpha_1, \dots, \alpha_n$, WEIGHTED-BINOMIAL-MEDIAN (WBM) asks for a median of the random variable

$$Z := \sum_{i=1}^n \alpha_i \text{Bernoulli}(x),$$

where the $\text{Bernoulli}(x)$ random variables are independent (and identically distributed).

This weighted binomial distribution (WBD) is comparable to the Poisson binomial distribution (PBD) in that they both generalize the binomial distribution. However the PBD is an unweighted sum of Bernoulli random variables with distinct probabilities x_i , while the WBD is a sum of Bernoulli random variables with a common x but distinct integer weights.

Lemma 6. *WBM is #P-Hard.*

The proof of the lemma is relegated to the appendix. We now give the main reduction.

Proof of Theorem 3. We reduce from WBM. If $k = 1$ then the reduction is immediate: if each of the P_i is a scaled down copy of $\alpha_i \text{Bernoulli}(x)$, then finding the optimal report for the random variable $\sum_i P_i$ amounts to finding the (scaled down) median of $\sum_i \alpha_i \text{Bernoulli}(x)$.

More generally, given an instance of WBM described by $(x, \alpha_1, \dots, \alpha_n)$, we will construct an instance of MAE-ESTIMATOR for any $k \geq 2$ for which determining optimal partitions and reporting scheme will solve our instance of WBM.

Our P_i will be discrete distributions given by

$$\Pr \left[p_i = \frac{1}{2k} \right] = \frac{1-x}{k} \quad (5)$$

$$\Pr \left[p_i = \frac{1 + \delta \frac{\alpha_i}{\sum_t \alpha_t}}{2k} \right] = \frac{x}{k} \quad (6)$$

$$\Pr \left[p_i = \frac{2j-1}{2k} \right] = \frac{1}{k} \quad \text{for } j = 2, \dots, k. \quad (7)$$

We will choose δ small enough such that the optimal partition of each of the P_i necessarily groups the atoms described in Equation (5) and Equation (6) together, and gives each of the atoms of Equation (7) its own interval in the partition. To find such a δ , first consider the “good” case when the partitions are of this form. In this case, there are k^n total boxes, each with weight $1/k^n$. Within each box C , the distribution of ℓ_1 norms has range upper bounded by $\delta/(2k)$. Within each C , the range of this distribution is an upper bound on the ℓ_1 distance between any atom in C and the optimal report for C . Therefore a loose upper bound on total MAE is

$$\sum_{c \in [k]^n} P(C_c) \frac{\delta}{2k} = \frac{\delta}{2k}. \quad (8)$$

On the other hand, consider the “bad” case when at least one of the partitions groups either two of the Equation (7)

atoms together or the Equation (6) atom together with at least one of the Equation (7) atoms. Assume without loss of generality that the $i = 1$ partitioning is “bad”. We will focus on the case when an Equation (6) and at least one Equation (7) atom are grouped together (because it is an interval, necessarily $j = 2$ is included), since in the best case it is the least costly scenario. Because of the product structure of the boxes induced by the partitions, for every pair of vectors u and u' in the support of P of the form

$$u = \left(\frac{1 + \delta \frac{\alpha_i}{\sum_j \alpha_j}}{2k}, u^- \right) \quad u' = \left(\frac{3}{2k}, u^- \right),$$

where $u^- \sim \prod_{j=2}^n P_j$, necessarily u and u' are contained in the same box. Therefore among each pair of u and u' , at least $M_{u^-} = \frac{\min\{x, 1-x\}}{k} \prod_{j=2}^n P_j(u_j^-)$ mass must travel $\|u'\|_1 - \|u\|_1$ to the estimate for their shared box, which yields a lower bound on the error of

$$\sum_{u^-} \left(\frac{1-\delta}{k} M_{u^-} \right) = \frac{(1-\delta) \min\{x, 1-x\}}{k^2}. \quad (9)$$

By Equations (8) and (9), choosing a $\delta < \frac{\min\{x, 1-x\}}{k}$ guarantees that the optimal partitioning for our instance is the “good” partitioning, and so all of the Equation (5) and Equation (6) atoms appear in the same box $C^* = \prod_i B_{i,1}$.

Recall that by Lemma 2, the MAE-minimizing estimate for a fixed box C is a median of the distribution of ℓ_1 norms of the vectors $u \in C$ according to P . Therefore MAE-ESTIMATOR finds some MAE-optimal report r^* for the box C^* , which by Equation (5) and Equation (6) implies that $\frac{r^* - n/2k}{\delta}$ is a median of $\sum_i \alpha_i \text{Bernoulli}(x)$, solving the given instance of WBM. \square

5 Discussion

Although we have assumed throughout that $p_i \in [0, 1]$, the results in Section 3 can be generalized to general bounded p_i . Similarly, the optimal strategy described in Section 4.1 holds for prior distributions over unbounded $p_i \in \mathbb{R}$.

There are also promising avenues for future work to extend our results to richer settings. For instance, in Section 4 we assume that the prior distributions P_i are given to us. However, it could be interesting to consider a setting in which the P_i are initially unknown but gradually discovered over rounds of questions; i.e., a learning setting where the elicitation scheme learns the prior distributions P_i in the course of accurately estimating the mean of the p_i .

Moreover, while we have focused on estimating $\|p\|_1$, one may ask if it is possible to estimate other functions of p . For instance, can the median of the p_i be efficiently and accurately estimated in this multiple-choice question setting?

Future work may also explore how our results carry over to settings in which the error metric is asymmetric: For instance, it may be more costly to underestimate than overestimate the size of a matriculating class due to space and resource constraints, and an optimal estimator would take this cost asymmetry into account.

References

- Acharya, J.; Diakonikolas, I.; Hegde, C.; Li, J. Z.; and Schmidt, L. 2015. Fast and near-optimal algorithms for approximating distributions by histograms. In *Proceedings of the 34th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*, 249–263.
- Babcock, B.; Chaudhuri, S.; and Das, G. 2003. Dynamic sample selection for approximate query processing. In *Proceedings of the 13th ACM SIGMOD/PODS International Conference on Management of Data*, 539–550.
- Berg, J.; Forsythe, R.; Nelson, F.; and Rietz, T. 2008. Results from a dozen years of election futures markets research. *Handbook of experimental economics results* 1(1):742–751.
- Boutilier, C. 2002. A pomdp formulation of preference elicitation problems. In *Proceedings of the 18th AAAI Conference on Artificial Intelligence (AAAI)*, 239–246.
- Brier, G. W. 1950. Verification of forecasts expressed in terms of probability. *Monthly weather review* 78(1):1–3.
- Chaudhuri, S.; Das, G.; Datar, M.; Motwani, R.; and Narasayya, V. 2001. Overcoming limitations of sampling for aggregation queries. In *Proceedings of the 17th IEEE International Conference on Data Engineering (ICDE)*, 534–542.
- Chaudhuri, S.; Das, G.; and Narasayya, V. 2007. Optimized stratified sampling for approximate query processing. *ACM Transactions on Database Systems (TODS)* 32(2): article 9.
- Chen, Y., and Kash, I. A. 2011. Information elicitation for decision making. In *Proceedings of the 10th International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*, 175–182.
- Garey, M. R., and Johnson, D. S. 1979. *Computers and Intractability: a Guide to the Theory of NP-Completeness*. W. H. Freeman and Company.
- Gneiting, T., and Raftery, A. E. 2007. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association* 102(477):359–378.
- Good, I. J. 1952. Rational decisions. *Journal of the Royal Statistical Society* 14(1):107–114.
- Gray, R., and Neuhoff, D. 1998. Quantization. *IEEE Transactions on Information Theory* 44(6):2325–2383.
- Hoeffding, W. 1963. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association* 58(301):13–30.
- Jagadish, H. V.; Koudas, N.; Muthukrishnan, S.; Poosala, V.; Sevcik, K. C.; and Suel, T. 1998. Optimal histograms with quality guarantees. In *Proceedings of the 24th International Conference on Very Large Data Bases (VLDB)*, 24–27.
- Jégou, H.; Douze, M.; and Schmid, C. 2011. Product quantization for nearest neighbor search. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 33(1):117–128.
- Jensen, R. E. 1969. A dynamic programming algorithm for cluster analysis. *Operations Research* 17(6):1034–1057.
- Lambert, N., and Shoham, Y. 2009. Eliciting truthful answers to multiple-choice questions. In *Proceedings of the 10th ACM Conference on Economics and Computation (EC)*, 109–118.
- Lloyd, S. 1982. Least squares quantization in PCM. *IEEE Transactions on Information Theory* 28(2):129–137.
- Radanovic, G., and Faltings, B. 2014. Incentives for truthful information elicitation of continuous signals. In *Proceedings of the 28th AAAI Conference on Artificial Intelligence (AAAI)*, 770–776.
- Soloviev, M., and Halpern, J. Y. 2018. Information acquisition under resource limitations in a noisy environment. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence (AAAI)*, 6443–6450.
- Waggoner, B., and Chen, Y. 2014. Output agreement mechanisms and common knowledge. In *Proceedings of the 2nd AAAI Conference on Human Computation and Crowdsourcing (HCOMP)*, 220–226.
- Winkler, R. L. 1969. Scoring rules and the evaluation of probability assessors. *Journal of the American Statistical Association* 64(327):1073–1078.
- Wolfers, J., and Zitzewitz, E. 2004. Prediction markets. *Journal of economic perspectives* 18(2):107–126.
- Yao, A. C. 1977. Probabilistic computations: Towards a unified measure of complexity. In *Proceedings of the 17th Symposium on Foundations of Computer Science (FOCS)*, 222–227.
- Zohar, A., and Rosenschein, J. S. 2008. Mechanisms for information elicitation. *Artificial Intelligence* 172(16–17):1917–1939.

A Proof of Lemma 6

In order to show to show that WBM is $\#\mathcal{P}$ -complete, we will reduce from the counting version of the knapsack problem, which is known to be $\#\mathcal{P}$ -complete (Garey and Johnson 1979):

Definition 1 ($\#\text{KNAPSACK}$). Given a list of nonnegative integer weights w_1, \dots, w_n and an integer capacity W , $\#\text{KNAPSACK}$ asks how many sets $S \subseteq [n]$ exist such that $\sum_{i \in S} w_i \leq W$.

And we will make use of a slight variant of counting knapsack:

Definition 2 ($\text{K}\#\text{KNAPSACK}$). Given an integer k , a list of nonnegative integer weights w_1, \dots, w_n an integer capacity W , and an integer threshold N , $\text{K}\#\text{KNAPSACK}$ finds $|S|$, where

$$S := \{S \subseteq [n] : \sum_{i \in S} w_i \leq W \text{ and } |S| = k\}.$$

Lemma 7. $\text{K}\#\text{KNAPSACK}$ is $\#\mathcal{P}$ -Complete.

Proof. There is an easy reduction from $\#\text{KNAPSACK}$: given an instance of $\#\text{KNAPSACK}$, simply query $\text{K}\#\text{KNAPSACK}$ for all values of k and return the sum of the answers. \square

Proof of Lemma 6. We begin by arguing that WBM may be assumed to return the largest possible median. This is because, for an instance of WBM given by $(x, \alpha_1, \dots, \alpha_n)$,

we may instead take a perturbed probability $\bar{x} = x + \gamma$. By choosing γ small enough, we can ensure that the median \bar{m} of $\bar{Z} := \sum_i \alpha_i \text{Bernoulli}(\bar{x})$ is a median of Z , but that it is the largest possible such median. Informally, we may tweak x gently enough that we preserve the median but break any median ties.

Formally, let F_Z be the cumulative density function (CDF) of Z . Since Z is a distribution comprised solely of atoms of weight $x^k(1-x)^{n-k}$ for $k \in [n]$, it suffices to find some perturbation γ for which

$$F_Z(m) - F_{\bar{Z}}(m) < a,$$

where m is a median of Z and a is a lower bound on the size of an atom in both Z and \bar{Z} . To show that we may choose such an a , note we may assume that $(1-x)^n \leq 1/2$, since otherwise the largest possible m is 0, and similarly that $\bar{x}^n \leq 1/2$, since otherwise we may easily check if the largest possible m is $\sum_{i \in [n]} \alpha_i$. Among all Z for which $x^n \leq 1/2$ and $(1-x)^n \leq 1/2$, the smallest possible atom is of size $\frac{1}{2}(2^{1/n} - 1)^n$, and so $a := 1/n^n$ is a lower bound on the atom size in Z for any value of x that concerns us.

Since Z is atomic, we then have that

$$F_Z(y) = \sum_{z \leq y} \Pr[Z = z] \quad (10)$$

$$= \sum_{S \subseteq [n]} x^{|S|} (1-x)^{n-|S|} \mathbb{1}_{\{\sum_{i \in S} w_i \leq y\}} \quad (11)$$

and so

$$\frac{\partial F_Z(y)}{\partial x} \leq \sum_{S \subseteq [n]} \frac{\partial}{\partial x} x^{|S|} (1-x)^{n-|S|} \leq n2^n. \quad (12)$$

Therefore taking $\gamma = \frac{a}{n2^n}$ will suffice, and $\bar{x} = x + \gamma$ will have a binary representation which is polynomial in the number of input bits.

We now reduce from K\#KNAPSACK . Given an instance of K\#KNAPSACK described by (k, w_1, \dots, w_n, W) , let $\Gamma := \langle k \rangle + \sum_i \langle w_i \rangle + \langle W \rangle$ be the length of the binary representation of these integers. For each i , let

$$\alpha_i := G + w_i,$$

where $G := (n+1) \sum_i w_i$. If $Z = \sum_i \alpha_i \text{Bernoulli}(x)$ for some rational $x \in [0, 1]$, then since the w_i are positive, the support of Z is clustered to the left of the integers $0, G, \dots, nG$. Specifically, we have by Equation (11) that

$$\begin{aligned} F_Z(Gk) &= \sum_{S \subseteq [n]} x^{|S|} (1-x)^{n-|S|} \mathbb{1}_{\{\sum_{i \in S} w_i \leq Gk\}} \\ &= \sum_{j=0}^{k-1} \binom{n}{j} x^j (1-x)^{n-j}, \end{aligned}$$

and so $F_Z(Gk)$ can be computed in time polynomial in $\Gamma + \langle x \rangle$.

Next, with k given, consider a binary search over (rational) x which searches for the largest possible x for which $m \leq Gk + W$. Once the binary search is far enough along

and the change in x is sufficiently small, $F_Z(m)$ approaches $1/2$ and the remaining change possible in $F_Z(m)$ will be small with respect to the atomic lower bound a . We may terminate our search, say, when $F_Z(m) \in [1/2, 1/2 + a/10]$. At this point m is the largest value of size at most $Gk + W$ in the support of Z , and so by this maximality of $m \leq Gk + W$,

$$\begin{aligned} F_Z(m) &= \sum_{S \subseteq [n]} x^{|S|} (1-x)^{n-|S|} \mathbb{1}_{\{\sum_{i \in S} w_i \leq m\}} \\ &= \sum_{j=0}^{k-1} \binom{n}{j} x^j (1-x)^{n-j} + |\mathcal{S}_k| x^k (1-x)^{n-k}. \end{aligned}$$

At this point a is much smaller than the other terms, and we may solve for $|\mathcal{S}_k|$, round, and solve K\#KNAPSACK :

$$|\mathcal{S}_k| \in \frac{1/2 \pm a/10 - \sum_{j=0}^{k-1} \binom{n}{j} x^j (1-x)^{n-j}}{x^k (1-x)^{n-k}}$$

It remains only to justify that this binary search for x terminates sufficiently quickly. By Equation (12) in order to guarantee that $F_Z(m)$ is within $a/10$ of $1/2$ it suffices to guarantee that the binary search step for x has size at most $\frac{a}{10n2^n}$. This requires $\log(10n^{n+1}2^n)$ steps, which is polynomial in n . \square