Optimal Scoring Rule Design

Yiling Chen*

Fang-Yi Yu[†]

Abstract

This paper introduces an optimization problem for proper scoring rule design. Consider a principal who wants to collect an agent's prediction about an unknown state. The agent can either report his prior prediction or access a costly signal and report the posterior prediction. Given a collection of possible distributions containing the agent's posterior prediction distribution, the principal's objective is to design a bounded scoring rule to maximize the agent's worst-case payoff increment between reporting his posterior prediction and reporting his prior prediction.

We study two settings of such optimization for proper scoring rules: static and asymptotic settings. In the static setting, where the agent can access one signal, we propose an efficient algorithm to compute an optimal scoring rule when the collection of distributions is finite. The agent can adaptively and indefinitely refine his prediction in the asymptotic setting. We first consider a sequence of collections of posterior distributions with vanishing covariance, which emulates general estimators with large samples, and show the optimality of the quadratic scoring rule. Then, when the agent's posterior distribution is a Beta-Bernoulli process, we find that the log scoring rule is optimal. We also prove the optimality of the log scoring rule over a smaller set of functions for categorical distributions with Dirichlet priors.

1 Introduction

Proper scoring rules are scoring functions that incentivize truthful information elicitation: an agent who has a subjective belief about an uncertain event maximizes his expected score by making a prediction according to his belief. When the agent can acquire costly information to refine his belief, which proper scoring rules maximally incentivize the agent's information acquisition? We formalize this question and show that popular quadratic and log scoring rules are optimal for information acquisition under different scenarios.

Suppose a principal wants to elicit a probabilistic prediction about an uncertain state W (e.g. the outcome of a future coin flip). An agent can report his prior P(W) (e.g. uniform on each side). Alternatively, the agent can obtain additional costly information S, say by tossing the coin several times, and then update his prediction to P(W|S=s). Any proper scoring rule can elicit the agent's prediction truthfully once its formed. But not all scoring rules incentivize the agent to gather costly information and improve his posterior prediction. In particular, the information gain — the expected difference of payment on posterior prediction and prior prediction — may be too small for the agent to obtain information. We propose an optimization framework to study the principal's problem of designing a proper scoring rule to incentivize the agent to acquire costly information, subject to a bounded payment constraint. The principal hopes to maximize the agent's information gain, with respect to the chosen scoring rule, for the worst information structure that the agent may have. This worst-case consideration is natural. For example, if the principal wants to design a scoring rule to incentivize her students to grade assignments with effort, the principal needs to ensure her scoring rule can incentivize all students to work hard when students have heterogeneous information structures.

We study two settings of this scoring rule design problem: static and asymptotic settings. In the static setting, the agent can access one signal. The principal's problem (Problem 1) is to design a proper scoring rule that 1) maximizes the information gain of the worst information structure in a known set \mathcal{P} that the agent's information structure belongs to, and 2) is subjected to a bounded expected payment constraint. We require the expected score to be always bounded as otherwise, the problem is meaningless because we can

^{*}Harvard University yiling@seas.harvard.edu

[†]Harvard University fangyiyu@seas.harvard.edu

always scale the scoring rule. Our setting is similar to that of Hartline et al. [20], but they require the ex-post payment be bounded. See related work for more comparison. We also consider the reverse of Problem 1 (our Problem 2) in the static setting and ask, for any given proper scoring rule, whether there exists a set of information structures that make the scoring rule uniquely optimal for incentivizing information acquisition.

In the asymptotic setting (Problem 3), the agent can adaptively and indefinitely refine his prediction. For instance, suppose we want to design a scoring rule and use it as a market scoring rule in a prediction market to encourage a sequence of forecasters to exert effort and report their posterior prediction. Each forecaster's incentive depends on the previous forecaster's information, and they can adaptively decide whether to spend the effort. Thus we need to design a scoring rule that optimizes for a sequence of collections of information structures where the n-th collection represents all possible information structures of any k-th forecaster with $k \leq n$. Another example is the large sample theory in statistics which studies estimators with an indefinitely growing sample size. In this context, Problem 3 asks what the optimal proper scoring rule is when the agent can adaptively get sample points to improve his prediction.

1.1 Our Technique and Contributions

Because any proper scoring rule can be specified as a convex function of prediction and its sub-gradient [22, 29, 18], our design space of the optimization in Problem 1 and 3 can be converted to the set of bounded convex functions, and the information gain becomes the Jensen gap of the associated convex function. This gives us a geometric interpretation of our optimization problem: The information gain amount to how curved (an integral of second derivative) the associated convex function, H, is at the information structure. However, the double iterative integral of the second derivative of H is bounded, since H is bounded.

Section 4 studies the static setting. For Problem 1, we first consider the extreme case when the principal knows the agent's information structure exactly (i.e. \mathcal{P} is singleton), and show in Section 4.1 that the optimal convex function is an upside-down pyramid (Fig. 1). Then Section 4.2 considers the case that the collection of information structures is finite: both the number of information structures and the support of any information structure are finite. Theorem 4.4 shows we can solve the optimization problem efficiently and the optimal H is a convex piecewise linear function. Section 4.3 looks at Problem 2: Given any convex H, is there a collection of information structures for which H is (uniquely) optimal? We show a necessary and sufficient condition for a convex piecewise linear function to be uniquely optimal in Theorem 4.6. However, we show quadratic scoring rule cannot be uniquely optimal in Problem 1.

Section 5 considers the asymptotic setting which generalizes the static setting and evaluates a scoring rule relatively against the best scoring rule called the relative information gain (Problem 3). Section 5.1 tries to study general Bayesian estimation with large number of i.i.d. samples. Inspired by (Bayesian) Cramèr-Rao bound [17] which shows the Bayes estimator with n samples has variance of order $\Theta(1/n)$, Section 5.1 considers the sequence of collections of information structures with vanishing covariance. Theorem 5.1 shows that the quadratic scoring rule optimizes the relative information gain against any smooth functions. To prove this, we relate the information gain to the strong convexity of our convex function and show that the quadratic scoring rule is the "most strongly convex" function.

Finally, Section 5.2 studies the information structures that are in an exponential family with conjugate prior, Beta-Bernoulli and Dirichlet-categorical information structures, and shows that the log scoring rule is optimal. Specifically, Section 5.2.1 studies a sequence of collections of information structures that converges to the collection of all Beta-Bernoulli distributions and show the relative information gain of the log scoring rule is optimal against all smooth functions. Additionally, through simulation Fig. 3 shows that the optimal scoring rules for Beta-Bernoulli distribution with large enough sample size is close to log scoring rule. Section 5.2.2 extends this result to non-binary settings and show the log scoring rule optimizes the relative information gain against all functions that are smooth up to the boundary. To achieve these results, we first show the limit of information gain is an elliptic differential operator on the convex function. Then we use a maximum principle to show that the log scoring rule maximized the elliptical differential operator's value [11, 4].

1.2 Related Work

Our problem can be seen as purchasing prediction from strategic people. The work on this topic can be roughly divided into two categories according to whether agents can misreport their signal or prediction. Below we focus on the relationship of our work to the most relevant technical scholarship.

In the first category, to ensure agents reporting their signals or predictions truthfully, there are two settings according to whether money is used for incentive alignment. In the first setting, the analyst uses monetary payments to incentivize agents to reveal their data truthfully. The challenge is to ensure truthtelling gets the highest payments. Existing works verify agents' reports by using either an observable ground truth (proper scoring rules) or peers' reports (peer prediction). For the second setting, individuals' utilities directly depend on an inference or learning outcome (e.g. they want a regression line to be as close to their own data point as possible) and hence they have incentives to manipulate their reported data to influence the outcome. [10, 23, 27, 19, 8].

Our setting is inspired by Hartline et al. [20] where there is an observable ground truth to verify an agent's reports and we want to maximize agent's inventive under a budget constraint. However, they consider that the ex-post payment is bounded, while we ask the ex-ante payment be bounded. This different boundedness condition contributes to the distinction between our results and theirs. In particular, ex-post bounded payment excludes the log scoring rule which is arguably one of the most important proper scoring rule. Furthermore, we show the log scoring rule is indeed the optimal when agent's belief is in exponential family with conjugate prior. Additionally, our paper stresses more on unknown information structure, and theirs focuses more on known information structure. They show if the information structure is exactly known the optimal scoring rule has v-shaped in the one-dimensional case which is qualitatively similar to our Corollary 4.2. They also consider general information structures and show that the quadratic scoring rule is optimal up to a constant factor which also holds for any strongly convex function. In contrast, we show the quadratic scoring rule is uniquely optimal.

Neyman et al. [24] also study optimal scoring rule design problem, but is more related to sequential method [16]. Instead of imposing bounded payment conditions, their objective comprises both payment and accuracy. They consider a special case of Beta-Bernoulli information structures where the prior is uninformative and design scoring rule for an agent to sequentially acquire samples from a Bernoulli distribution.

In the second category, agents cannot misreport their signal or prediction. The problem of purchasing data from people has been investigated with different focuses, e.g. privacy concerns [14, 12, 15, 25, 9, 30], effort and cost of data providers [28, 5, 1, 7, 31, 6], and reward allocation [13, 3].

2 Problem Description

A principal wants to collect high quality prediction of an unknown state of the world W from an agent who can access a costly signal S. The space of the world Ω is finite with $|\Omega| = d$ and has a typical realization $\omega \in \Omega$. The signal space S is a measurable set with a typical realization $s \in S$. We refer to a joint distribution P on $\Omega \times S$ as an information structure. Because accessing the signal S is costly for the agent, the principal wants to design a scoring rule to incentivize the agent to spend the effort and report the posterior prediction instead of the prior prediction. However, the principal only knows a collection of information structures $P \subseteq \Delta(\Omega \times S)$ that the agent's information structure falls into.

Before stating our objective, we first simplify the notation by calibrating agent's signal. Specifically, let the agent's prediction after observing the signal be a random variable on $\Delta(\Omega)$, $X = X(s) := P(W \mid S = s)$ when S = s, and, thus, his prediction without the signal is $\mathbb{E} X = P(W)$. For the principal, the information structures $P \in \mathcal{P}$ affects only the random variable of prediction X. Hence, we will use P and X exchangeably for information structures in this paper, and let \mathcal{X} be the collection of random variables of prediction induced from information structures P in \mathcal{P} .

With the above definitions, the principal wants to maximize the difference of expected payment between reporting the prior prediction $\mathbb{E} X$ and reporting the posterior prediction X. Because a proper scoring rule $PS: \Omega \times \Delta(\Omega) \to \mathbb{R}$ can be associated with a convex function $H: \Delta(\Omega) \to \mathbb{R}$ so that

$$PS(\omega, x) = H(x) + \partial H(x) \cdot (\mathbf{1}_{\omega} - x) \tag{1}$$

for all $x \in \Delta(\Omega)$ and $\omega \in \Omega$ (See Section 3.1 for more properties of scoring rules), given an information structure X, the value of observing the signal under the proper scoring rule with H is

$$\mathfrak{J}_H(X) := \mathbb{E}_X[H(X)] - H(\mathbb{E} X).$$

We call $\mathfrak{J}_H(X)$ information gain of X on H which measures the difference of expected scores between agents who exert effort and those who do not.¹

Given Ω and \mathcal{X} , the principal chooses a convex function H which maximizes the worst information gain of \mathcal{X} which is called the information gain of \mathcal{X} on H,

$$\mathrm{Obj}_{\mathcal{X}}(H) := \inf_{X \in \mathcal{X}} \mathfrak{J}_H(X).$$

We require the scoring rules are bounded so that the ex-ante payment of truth-telling is always in [0, 1]. By a simple calculation, a bounded scoring rule has H in $\mathcal{H} := \{H : 0 \leq H(x) \leq 1, \forall x \in \Delta(\Omega)\}$.

Now we are ready to define our max-min optimization problem for the principal. First, Section 4 considers a max-min optimization problem given a collection of information structures.

Problem 1 (static). Given a set of information structures \mathcal{X} and Ω , find bounded H which maximizes the information gain of \mathcal{X}

$$\operatorname{Obj}_{\mathcal{X}}(H) = \max_{\tilde{H} \in \mathcal{H}} \operatorname{Obj}_{\mathcal{X}}(\tilde{H}).$$
 (2)

We further call $\operatorname{Opt}(\mathcal{X}) := \max_{H \in \mathcal{H}} \operatorname{Obj}_{\mathcal{X}}(H).^2$

We can also ask the converse of Problem 1: Given any convex H, is there a collection of information structures \mathcal{X} for which H is (uniquely) optimal in Problem 1? These problem can provide an incentive understanding why certain proper scoring rules, e.g., log scoring rule and quadratic scoring rule are commonly used.

We say a collection of information structures \mathcal{X} settles a convex function $H \in \mathcal{H}$ if H is optimal for \mathcal{X} . Given $U \subseteq \Delta_d$, we say \mathcal{X} uniquely settles H on U if for all $H^* \in \arg \max_{\tilde{H}} \operatorname{Obj}_{\mathcal{X}}(\tilde{H})$, $H^*(x) = H(x)$ for all $x \in U$, and uniquely settles H when $U = \Delta_d$. Equivalently H is the unique optimal for \mathcal{X} , $\{H\} = \arg \max_{\tilde{H}} \operatorname{Obj}_{\mathcal{X}}(\tilde{H})$.

Problem 2. Given $U \subseteq \Delta_d$ and a convex function H, find a collection of information structure \mathcal{X} so that \mathcal{X} uniquely settles H, i.e., H is the unique optimal for \mathcal{X} .

Note that Problem 2 requires the unique optimality. Otherwise, if a collection of information structures contains a constant information structure X where $X = \mathbb{E} X$, any H has zero information gain of X, and thus is optimal.

The agent may adaptively conduct a sequence of experiments and access a sequence of signals. Hence the principal may want the agent to access as many signals as possible, and needs to ensure the information gain of each additional signal is as large as possible.

Problem 3 (asymptotic). Let \mathcal{X}_0, \ldots be a sequence of collections of information structures which does not contain any constant information structure and $\bar{\mathcal{X}}_n := \bigcup_{k \leq n} \mathcal{X}_k$ for all n. The principal chooses H to maximize the *relative information gain* of the sequence \bar{X}_0, \ldots ,

$$\inf_{\tilde{H} \in \mathcal{H}} \lim_{n \to \infty} \frac{\mathrm{Obj}_{\bar{\mathcal{X}}_n}(H)}{\mathrm{Obj}_{\bar{\mathcal{X}}_n}(\tilde{H})}.$$
 (3)

As the example in the introduction, to predict an output of a biased coin, the agent can observe any number of samples from the coin. The principal has a sequence of collections of information structures $\mathcal{X}_0, \mathcal{X}_1, \ldots$ where each element in \mathcal{X}_n is a possible information structure of the (n+1)-th sample. Specifically, given an outcome $s^{(\leq n)} = (s^{(1)}, \ldots, s^{(n)}) \in \{0, 1\}^n$ of the first n signal $S^{(\leq n)}$, let $X_{s^{(\leq n)}}(s) = P(W \mid S^{(\leq n)} = s^{(\leq n)}, S^{(n+1)} = s)$ be the posterior on W with the (n+1)-th signal conditional on $S^{(\leq n)} = s^{(\leq n)}$ or posterior

 $^{{}^{1}\}mathfrak{J}_{X}(H)$ can be seen as the value of information for the agent [21], or Jensen's gap of the random variable X on H. [2]

 $^{^2}$ Eq. (2) has the maximum, because the space \mathcal{H} with the sup norm $\|\cdot\|_{\infty}$ is complete and for any \mathcal{X} the functional $Obj_{\mathcal{X}}$ is continuous on the norm space $(\mathcal{H}, \|\cdot\|_{\infty})$.

with (n+1)-th signal for short. Then $\mathcal{X}_n = \{X_{s^{(\leq n)}}, \forall s^{(\leq n)} \in \{0,1\}^n\}$. As the asymptotic theory in statistics studies estimators with indefinitely growing sample size, Problem 3 evaluates proper scoring rules when the agent may collect indefinitely data. The notion of relative information gain is particularly useful when the cost of each signal is small for all, or the number of signal is indefinite.

Equation (3), intuitively, measures the information gain of $\lim_{n\to\infty} \bar{\mathcal{X}}_n$ on H. Indeed Proposition 2.1 shows that Problem 3 contains Problem 1 as a special case when $\operatorname{Opt}(\lim_n \bar{X}_n) > 0$. Furthermore, using the ratio of information gains, Eq. (3) is meaningful even when $\operatorname{Opt}(\lim_n \bar{X}_n) = 0$.

Proposition 2.1. Given an increasing sequence \mathcal{X}_1, \ldots with $\mathcal{X} := \lim_{n \to \infty} \mathcal{X}_n$, if $\operatorname{Opt}(\mathcal{X}) > 0$, $H \in \mathcal{H}$ is optimal for \mathcal{X} for Problem 1 if and only if H is optimal for the sequence \mathcal{X}_1, \ldots for Problem 3.

3 Preliminary

Here we list some notations. Given positive integer d and n, let $[n] := \{1, \ldots, n\}$, $\Delta = \Delta_d = \Delta(\Omega) \subset \mathbb{R}^d$ is probability simplex over $\Omega = [d]$. The vertices of the simplex Δ_d are \hat{e}_j for $j \in [d]$ which is also the standard basis of \mathbb{R}^d . We set the interior and the boundary of Δ_d to be $\operatorname{int}(\Delta_d) := \{x \in \Delta_d : x_k > 0, \forall k \in [d]\}$, and $bd(\Delta) := \Delta \setminus \operatorname{int}(\Delta)$ respectively.

We say a function $h \in \mathcal{C}^{\infty}(U)$ is smooth on an open set U if h is all differentiable for all degrees of differentiation in U, and $h \in \mathcal{C}^{\infty}(\overline{U})$ is smooth on an close set \overline{U} if all differentiation of h is uniformly continuous in \overline{U} .

We further call all one vector $\mathbf{1} := (1, \dots, 1) \in \mathbb{R}^d$ and $c := \frac{1}{d}\mathbf{1}$. The 2-norm of a matrix A is $||A||_2 := \max_{v:||v||_2=1} ||Av||_2$.

3.1 Proper Scoring Rules

A scoring rule for a random variable $W \in \Omega$ is a function $PS : \Omega \times \Delta(\Omega) \to \mathbb{R}$ where $PS(\omega, \hat{x})$ is the score assigned to a prediction $\hat{x} \in \Delta(\Omega)$ when $W = \omega$. The scoring rule is (strict) proper if for all random variable W on Ω with distribution x, setting $\hat{x} = x$ (uniquely) maximizes the expected score $\mathbb{E}_{W \sim x} PS(W, \hat{x})$. In other words, if W is distributed according to x, then truthfully reporting x can maximize the expected score.

Theorem 3.1 (Savage representation [22, 29, 18]). For every (strict) proper scoring rule PS, there exists a (strictly) convex function $H: \Delta(\Omega) \to \mathbb{R}$ so that for all $\omega \in \Omega$ and $x \in \Delta(\Omega)$

$$PS(\omega, x) = H(x) + \partial H(x) \cdot (\mathbf{1}_{\omega} - x)$$

where $\partial H(x)$ is a sub-gradient of H at x and $\mathbf{1}_{\omega} \in \Delta(\Omega)$ is the indicator function with $\mathbf{1}_{\omega}(w') = 1$ if $w' = \omega$ and 0 otherwise.

Conversely, for every (strictly) convex function $H:\Delta(\Omega)\to\mathbb{R}$, there exists a (strict) proper scoring rule such that the above condition hold.

We list some common proper scoring rules with associated convex functions which are scaled to be in \mathcal{H} .

quadratic scoring rule
$$Q(x) = \frac{d}{d-1} \sum_{k=1}^{d} (x_k^2 - 1/d)^2$$

spherical scoring rule $H_s(x) = \frac{\sqrt{d}}{\sqrt{d}-1} \sqrt{\sum x_k^2} - \frac{1}{\sqrt{d}-1}$

log scoring rule
$$H_{\ln}(x) = \frac{1}{\ln d} \sum_k x_k \ln x_k + 1$$

Note that the convex functions Q and H_s for quadratic scoring rule and spherical scoring rule respectively are smooth on the closed set Δ_d , but H_{ln} for log scoring rule is only smooth on the open set $int(\Delta_d)$.

When the state space Ω is binary, the associated convex function is one-dimensional. For every proper scoring rule PS, there exists a convex function $H:[0,1]\to\mathbb{R}$ so that for all $x\in[0,1]$ and binary event $\omega\in\{0,1\}$

$$PS(\omega, x) = \begin{cases} H(x) + \partial H(x) \cdot (1 - x) & \text{when } \omega = 1, \\ H(x) - \partial H(x) \cdot x & \text{when } \omega = 0. \end{cases}$$
(4)

3.2 Convex Functions

In Sections 5.1 and 5.2, we will restrict our functions to be smooth and utilize their analytic properties (k-th order derivative). Though the domain of our convex functions is a probability simplex Δ_d , which is a submanifold in \mathbb{R}^d , we can extend any smooth function h on Δ_d to a smooth function on $\mathbb{R}^d_{\geq 0} := \{x \in \mathbb{R}^d : x_k \geq 0, \forall k \in [d]\},^3$, and use standard notation to discuss the k-th derivative of h, instead of using complex coordinate charts.

We first define the notion of strongly convex on $int(\Delta)$.

Definition 3.2. H is α -strongly convex on $\operatorname{int}(\Delta)$, if for all $x, y \in \operatorname{int}(\Delta)$

$$H(y) \ge H(x) + \nabla H(x)(y - x) + \frac{\alpha}{2} ||y - x||^2.$$
 (5)

Note that the pair x and y are restricted in $int(\Delta)$, so the equivalent condition on the Hessian of H is also restricted in the tangent space of Δ . Formally,

Lemma 3.3 (Hessian). Suppose $H \in \mathcal{H} \cap \mathcal{C}^{\infty}(\operatorname{int}(\Delta))$. H is α -strongly convex on $\operatorname{int}(\Delta)$ if and only if for all $x \in \operatorname{int}(\Delta)$ and $v \in \mathbb{R}^d$ with $\mathbf{1}^{\top}v = 0$, $v^{\top}\nabla^2 H(x)v \geq \alpha \|v\|^2$.

3.3 PDE and Maximum Principle

We will associate the objective in Problem 3 with elliptic differential operators and exploit the maximum principles of our differential operator.

Definition 3.4 (Evans [11]). Let $U \subset \mathbb{R}^d$ be an open and bounded set. \mathcal{L} is a second order differential operator with $A: U \to \mathbb{R}^{d \times d}$ on a field $u: U \to \mathbb{R}$ if for all $x \in U$

$$\mathcal{L}u(x) = \sum_{i,j} A_{ij}(x)\partial_{ij}^2 u.$$

We say \mathcal{L} is *elliptic* if $A(x) = (A_{ij}(x))$ is positive definite for all $x \in U$, and *uniformly elliptic* if there exists $\rho > 0$ such that $A(x) - \rho \mathbb{I}$ is still positive definition for all $x \in U$.

The following theorem shows the maximum of a (super-)solution of a second order elliptic differential equation can only have its maximum points on the boundary. We will strengthen this result to functions on a probability simplex (Lemma 5.16).

Theorem 3.5 (maximum principle of elliptic PDE [11]). Suppose that \mathcal{L} is uniformly elliptic and $\mathcal{L}u \geq 0$ for a function $u \in \mathcal{C}^2(U)$. Then the function u must attain its maximum on the boundary

$$\max_{cl(U)} u = \max_{bd(U)} u$$

Moreover, if u has an interior maximum, u is a constant.

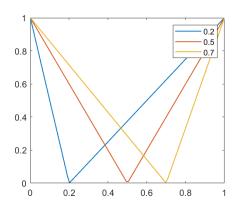
4 Static Setting

4.1 Singleton Information Structure

As a warm-up, let's consider the principal exactly knows the agent's information structure so that $\mathcal{X} = \{X\}$ is a singleton. We show the optimal H can be an upside down pyramid. (Fig. 1)

Theorem 4.1 (singleton). If $\mathcal{X} = \{X\}$ is singleton and the state space $\Omega = [d]$, there exists an optimal scoring rule associated with an upside down pyramid H^* such that the epigraph of H^* is the convex hull with vertices $(\mathbb{E}X,0)$, and $(\hat{e}_k,1)$ for all $k \in [d]$.

³Besides using Whitney's Expansion Theorem, for example, we can also define $(x_1, \ldots, x_d) \mapsto ||x|| \cdot h\left(\frac{x_1}{\sum_{k=1}^d x_k}, \ldots, \frac{x_d}{\sum_{k=1}^d x_k}\right)$ when $x \neq \mathbf{0}$, and 0 when $x = \mathbf{0}$. Furthermore, several common scoring rules are already defined on $\mathbb{R}^d_{>0}$, e.g., log scoring rule and quadratic scoring rule.



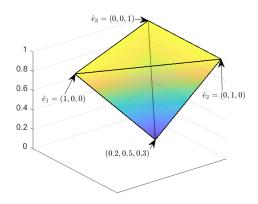


Figure 1: The left panel shows several examples of the optimal H in Corollary 4.2 which have $\mathbb{E} X = 0.2, 0.5$, and 0.7 respectively. The right panel shows the optimal H^* in Theorem 4.1 when d = 3, $\mathcal{X} = \{X\}$, and $\Pr[X = 1] = 0.2, \Pr[X = 2] = 0.5, \Pr[X = 3] = 0.3$.

Since the information gain $\mathfrak{J}_H(X)$ is an integration of $H(x) - H(\mathbb{E}X)$ over all $x \in \Delta_d$ for any H, we can prove Theorem 4.1 by a pointwise inequality, $H(x) - H(\mathbb{E}X) \leq H^*(x) - H^*(\mathbb{E}X)$ for all x.

Proof of Theorem 4.1. Let H^* be the upside down pyramid in Theorem 4.1, and $H \in \mathcal{H}$ be an arbitrary bounded convex function. Since $\mathcal{X} = \{X\}$ is singleton, it is sufficient to prove $\mathfrak{J}_H(X) \leq \mathfrak{J}_{H^*}(X)$.

Let $h_0 := H(\mathbb{E}X)$. If $h_0 = 1$, H is a constant function and $0 = \mathfrak{J}_H(X) \leq \mathfrak{J}_{H^*}(X)$ by the Jensen's inequality. Now we consider $h_0 < 1$. Let \tilde{H} be a convex piecewise linear function whose epigraph has vertices at $(\mathbb{E}X, h_0)$, and $(\hat{e}_k, 1)$ for all $k \in [d]$. First, because $H \in \mathcal{H}$, $H(\hat{e}_k) \leq \tilde{H}(\hat{e}_k) = 1$ for all $k \in [d]$, and $H(\mathbb{E}X) = \tilde{H}(\mathbb{E}X) = h_0$. Then, the epigraph of \tilde{H} is contained in the epigraph of H, so

$$H(x) - H(\mathbb{E}X) \le \tilde{H}(x) - \tilde{H}(\mathbb{E}X) \text{ for all } x \in \Delta_d.$$
 (6)

Second, because the epigraphs of H^* and \tilde{H} are both upside down pyramids, and the vertices are aligned, we can convert \tilde{H} to H^* through an affine transformation: $H(x) = \frac{\tilde{H}(x) - h_0}{1 - h_0}$ for all $x \in \Delta_d$. Therefore, for all $x \in \Delta_d$,

$$\tilde{H}(x) - \tilde{H}(\mathbb{E}X) = (1 - h_0) \left(H^*(x) - H^*(\mathbb{E}X) \right) \le H^*(x) - H^*(\mathbb{E}X), \tag{7}$$

because $0 \le h_0 < 1$, and both sides are non-negative. Combining Eqs. (6) and (7), we have $\mathfrak{J}_H(X) = \mathbb{E}_X [H(X) - H(\mathbb{E}X)] \le \mathbb{E}_X [H^*(X) - H^*(\mathbb{E}X)] = \mathfrak{J}_{H^*}(X)$, and complete the proof.

When the state space is binary $\Omega = \{0, 1\}$, by Eq. (4), the upside down pyramid H is v-shaped, and Theorem 4.1 yields the following corollary.

Corollary 4.2. If the state space $\Omega = \{0,1\}$ and $\mathcal{X} = \{X\}$, there exists an optimal scoring rule associated with a v-shape H such that

$$H(x) = \begin{cases} \frac{-1}{\mathbb{E}X} (x - \mathbb{E}X) & \text{if } x < \mathbb{E}X \\ \frac{1}{1 - \mathbb{E}X} (x - \mathbb{E}X) & \text{if } x \ge \mathbb{E}X \end{cases}$$

These results suggest the principal should choose H that is "curved" at the prior in order to incentivize the agent to derive the signal and move away from the prior. This intuition is useful for the later sections.

4.2 Finite Information Structures

In this section, we give a polynomial time algorithm that computes an optimal scoring rule when the collection of information structures is finite as defined below.

Definition 4.3. We call a collection of information structures \mathcal{X} finite if $|\mathcal{X}|$ is finite and all $X \in \mathcal{X}$ has a finite support $|\operatorname{supp}(X)| < \infty$.

When \mathcal{X} is finite, let $\overline{\operatorname{supp}}(X) := \operatorname{supp}(X) \cup \{\mathbb{E} X\} \subset \Delta(\Omega) \text{ for all } X \in \mathcal{X}, \text{ and } \overline{\operatorname{supp}}(\mathcal{X}) := \cup_{X \in \mathcal{X}} \overline{\operatorname{supp}}(X).$

The notion of finite collection of information structures is natural when there is a finite number of heterogeneous agents with finite set of information.

Theorem 4.4. If the state space $\Omega = [d]$, and \mathcal{X} is finite with $|\overline{\text{supp}}(\mathcal{X})| = m$, there exists an algorithm that computes an optimal bounded proper scoring rule and the running time is polynomial in d and m.

The main idea is that when \mathcal{X} is finite $\mathrm{Obj}_X(H)$ in Eq. (2) only depends on the evaluations of H in $\overline{\mathrm{supp}}(\mathcal{X})$. Thus, instead of searching for all possible bounded scoring rules, we can reduce the dimension of Problem 1 and use a linear programming whose variables contain the evaluations of H in $\overline{\mathrm{supp}}(\mathcal{X})$ and add linear constraints to ensure those evaluation can be extended to a convex function. This observation allows us the solve the problem in weakly polynomial time which is polynomial in d and m but may not be polynomial in the representation size.

Proof of Theorem 4.4. The idea is to construct a linear programming whose variables contain the evaluations of H in $\overline{\text{supp}}(\mathcal{X})$. To formulate this, we introduce some notations. Given $|\mathcal{X}| = n$, we set $\mathcal{X} = \{X_i : i = 1, \ldots, n\}$. For each $i \in [n]$, let the support of X_i be $\text{supp}(X_i) = \{x_{i,j} : j \in [m_i]\}$ with size $|\text{supp}(X_i)| = m_i$. Additionally, let the expectation be $x_{i,0} = \mathbb{E}[X_i]$, and $\Pr(X_i = x_{i,j}) = p_{i,j}$. Hence $\overline{\text{supp}}(X_i) = \{x_{i,j} : j = 0, \ldots, m_i\}$ and $\overline{\text{supp}}(\mathcal{X}) = \{x_{i,j} : i \in [n], j = 0, \ldots, m_i\}$. We further use $\mathcal{A} = \{(i,j) : i \in [n], j = 0, \ldots, m_i\}$ to denote the set of indices. Finally, we set the vertices of the probability simplex Δ_d be $x_k = \hat{e}_k$ for all $k \in [d]$, and $\bar{\mathcal{A}} := \mathcal{A} \cup [d]$. To simplify the notations, we assume $\overline{\text{supp}}(\mathcal{X})$ does not contain any vertex of Δ_d \hat{e}_k for $k = 1, \ldots, d$ and for all distinct α, α' in $\mathcal{A}, x_{\alpha} \neq x_{\alpha'}$.

Note that the objective value only depends on a finite number of values. Specifically, given $H(x_{\alpha}) = h_{\alpha}$ for any $\alpha \in \mathcal{A}$, the objective, Eq. (2), is

$$Obj_{\mathcal{X}}(H) = \min_{i \in [n]} \sum_{k=1}^{m_i} p_{i,j} h_{i,j} - h_{i,0}.$$
 (8)

Thus, we can first decide h_{α} to maximize Eq. (8), and "connect" those points (x_{α}, h_{α}) to construct a piecewise linear function. To ensure the resulting function is convex, we further require there exists a supporting hyperplane for each (x_{α}, h_{α}) — for each α there exists $g_{\alpha} \in \mathbb{R}^d$ such that $h_{\alpha'} \geq h_{\alpha} + g_{\alpha}^{\top}(x_{\alpha'} - x_{\alpha})$ for all $\alpha' \neq \alpha$.

In summary, we set the convex function to be

$$H(x) = \max \left\{ \max_{\alpha \in \bar{\mathcal{A}}} h_{\alpha} + g_{\alpha}(x - x_{\alpha}), \min_{\alpha} h_{\alpha} \right\}, \tag{9}$$

and the collection of h_{α} and g_{α} is a solution of the following linear programming,

$$\max \qquad \min_{i \in [n]} \sum_{k=1}^{m_i} p_{i,j} h_{i,j} - h_{i,0},
\text{subject to} \quad h_{\alpha} \in [0,1], \qquad \forall \alpha \in \bar{\mathcal{A}},
\quad h_{\alpha'} \ge h_{\alpha} + g_{\alpha}^{\top} (x_{\alpha'} - x_{\alpha}), \quad \forall \alpha, \alpha' \in \bar{\mathcal{A}}.$$
(10)

The above linear programming has $(1+d)|\bar{\mathcal{A}}| = O(d(m+d))$ variables and $|\bar{\mathcal{A}}| + |\bar{\mathcal{A}}|^2 = O((m+d)^2)$ constraints, so we can solve it in polynomial time with respect to m and d.

Now we need to show H is 1) convex, 2) bounded in [0, 1], and 3) optimal. It is easy to see for all $\alpha \in \bar{\mathcal{A}}$,

$$H(x_{\alpha}) = h_{\alpha},\tag{11}$$

⁴Otherwise, we just need to add some equality constraints. For instance, if $x_{\alpha} = x_{\alpha'}$, we need to set $h_{\alpha} = h_{\alpha'}$ as a constraint.

because

$$H(x_{\alpha}) = \max \left\{ h_{\alpha}, \max_{\alpha' \neq \alpha} h_{\alpha'} + g_{\alpha'}(x_{\alpha} - x_{\alpha'}), \min_{\alpha'} h_{\alpha'} \right\}$$

$$= \max \left\{ h_{\alpha}, \max_{\alpha' \neq \alpha} h_{\alpha'} + g_{\alpha'}(x_{\alpha} - x_{\alpha'}) \right\}$$

$$= h_{\alpha}$$
(by the constraints in Eq. (10))

First because H is the maximum of a collection of linear functions, H is convex. Second, for the lower bound, by the constraints in Eq. (10) $\min h_{\alpha} \geq 0$ so $0 \leq \min h_{\alpha} \leq H(x)$ due to Eq. (9). For the upper bound, because H is convex, for all $x \in \Delta_d$, $H(x) \leq \max_k \{H(\hat{e}_k)\} = \max_k \{h_k\} \leq 1$ by Eqs. (10) and (11). Finally, for any bounded convex function $\tilde{H} \in \mathcal{H}$, we set $\tilde{h}_{\alpha} = \tilde{H}(x_{\alpha})$ for $\alpha \in \bar{\mathcal{A}}$. At each x_{α} we can find a vector \tilde{g}_{α} such that $\tilde{H}(x) \geq \tilde{H}(x_{\alpha}) + \tilde{g}_{\alpha}^{\top}(x - x_{\alpha})$ for all $x \in \Delta_d$. Since \tilde{H} is convex and in \mathcal{H} , the collection of \tilde{h}_{α} and \tilde{g}_{α} is a feasible solution to Eq. (10), and $\mathrm{Obj}_{\mathcal{X}}(\tilde{H}) \leq \mathrm{Obj}_{\mathcal{X}}(H)$.

Note that if \mathcal{X} contains a constant information structure X, the resulting linear programming (10) will output an arbitrary piecewise linear convex function, because the objective value is always zero (Footnote 4).

4.3 Unique Optimality of an Convex Function

In this section we ask given any convex H, is there a collection of information structures \mathcal{X} that uniquely settles H? (Problem 2)

First Proposition 4.5 shows a simple necessary condition of H to be uniquely settled for some \mathcal{X} . For convex piecewise linear function, Theorem 4.6 shows such necessary condition is also sufficient, and we can construct a finite collection of information structures (Definition 4.3) to uniquely settle those convex piecewise linear functions. Similarly, Proposition 4.7 shows any convex H satisfying the necessary condition can be uniquely settled on any closed set U that does not contain any minimum or maximum point of H. However, Proposition 4.8 shows no collection of information structures uniquely settles a quadratic function which contrasts the unique optimality of quadratic scoring rules in an asymptotic setting in Section 5.1.

Proposition 4.5 (A necessary condition). For any bounded convex function H on Δ_d , there exists a collection of information structures \mathcal{X}_H with $\mathrm{Opt}(\mathcal{X}_H) > 0$ and settles H, only if

$$\min_{x} H(x) = 0, \text{ and } H(\hat{e}_{j}) = 1, \forall j \in [d].$$
(12)

The formal proof is in Appendix B. Intuitively, we can show that if one of the condition does not hold, we can construct an other convex function that has a larger or the same information gain of any information structure.

Convex piecewise linear functions Now, we consider sufficient conditions for piecewise linear functions H. Theorem 4.6 shows the necessary conditions Eq. (12) is also sufficient for any convex piecewise linear functions.

Theorem 4.6. If H is a convex piecewise linear function satisfying Eq. (12), there exists a finite collection of information structures \mathcal{X}_H that uniquely settles H.

Our proof is constructive: Given a point $\theta \in \Delta_d$, we can create information structures to upper bound and lower bound the evaluation of an optimal convex function. The proofs and details are in Appendix B.

General convex functions We can also use the same idea on a general convex function. Proposition 4.7 shows that we can construct a collection of information structures such that the optimal H^* agrees with H with all except $\theta \in \Delta_d$ with extreme large or small $H(\theta)$.

⁵Specifically, we can construct \tilde{g}_{α} by finding a support hyperplane to the epigraph of \tilde{H} at (x_{α}, h_{α}) , and the vector \tilde{g}_{α} is called subgradient.

Proposition 4.7. If H is a convex function satisfying Eq. (12), for any $\delta > 0$ there exists a collection of information structures $\mathcal{X}_{H,\delta}$ which uniquely settles H on $U_{H,\delta} := \{\theta \in \Delta_d : \delta < H(\theta) < 1 - \delta\}$.

Moreover, for all closed set $U \subset \Delta_d$ that does not contain any minimum or maximum point of H, there exists $\mathcal{X}_{H,U}$ that uniquely settles H on U.

However, Proposition 4.8 shows this limit is inevitable and no collection of information structures uniquely settles a quadratic scoring rule which satisfies Eq. (12). Finally, it's not hard to extend this example to show that no collection of information structures can settle any strictly convex function.

Proposition 4.8. For any collection of information structures \mathcal{X} which settles the quadratic scoring rule on binary state space $Q_2(x) = 4(x - \frac{1}{2})^2$, there exists another optimal H for \mathcal{X} but $H \neq Q_2$.

Proof of Proposition 4.8. Let $\delta := \operatorname{Opt}(\mathcal{X})$. If $\delta = 0$, any convex function is optimal and we are done. If $\delta > 0$, we define $x_{\delta} := \frac{1 - \sqrt{1 - \delta}}{2}$

$$H(x) = \begin{cases} Q_2(x) & \text{if } x \ge x_{\delta} \\ 1 - \frac{2\delta}{1 - \sqrt{1 - \delta}} x & \text{otherwise} \end{cases}$$

It's easy to check H is convex, $H(x) > Q_2(x)$ for all $0 < x < x_{\delta}$ and $H(x) = Q_2(x)$ otherwise.

Now we show a stronger result, $\mathfrak{J}_X(H) \geq \mathfrak{J}_X(Q_2)$ for all $X \in \mathcal{X}$. Suppose there exists $X \in \mathcal{X}$ with $\mathfrak{J}_X(H) < \mathfrak{J}_X(Q_2)$. If $\mathbb{E} X \geq x_\delta$, $\mathfrak{J}_X(H) = \mathbb{E} H(X) - H(\mathbb{E} X) = \mathbb{E} H(X) - Q_2(\mathbb{E} X) \geq \mathbb{E} Q_2(X) - Q_2(\mathbb{E} X) = \mathfrak{J}_X(Q_2)$ that contradicts to $\mathfrak{J}_X(H) < \mathfrak{J}_X(Q_2)$. Then if $\mathbb{E} X < x_\delta$, $\mathfrak{J}_X(Q_2) = \mathbb{E} Q_2(X) - Q_2(\mathbb{E} X) \leq 1 - Q_2(\mathbb{E} X) < \delta$ that contradicts to $X \in \mathcal{X}$ because $\mathrm{Opt}(\mathcal{X}) = \delta$.

5 Asymptotic Setting

Now we study optimal incentive scoring rules for eliciting estimation from Bayesian agents when the number of samples n is large. Because as the number of samples increases, the mean squared error (variance) is vanishing for any consistent estimator, $\operatorname{Opt}(\lim_{n\to\infty} \mathcal{X}_n) = 0$, we use the relative information gain in Problem 3 to evaluate a proper scoring rule.

Section 5.1 considers a sequence of collection of information structures with vanishing covariance, and shows the quadratic scoring rule is uniquely optimal. In Section 5.2, the data are generated from exponential families with conjugate prior (Beta-Bernoulli and Dirichlet-Categorical process), and the log-scoring rule is uniquely optimal. Though both Sections 5.1 and 5.2 consider eliciting Bayes estimator with large samples, Section 5.1 models general prior and observation which can be continuous, while the setting in Section 5.2 is more restricted, e.g., observation is in [d].

5.1 Vanishing Covariance Information Structures

Consider the agent can access n i.i.d. samples to estimate the distribution of the ground state W. What is a natural sequence of collections of information structures to model this process? By Bayesian Cramér-Rao lower bound [17], the variance (covariance) of the posterior distribution $P(W|S^{(1)}, \ldots S^{(n)})$ is of order $\Theta(1/n)$ for large n. Thus, the variance of the posterior with n-th signal given the first n-1 samples is of order $\Theta(1/n^2)$. This motivates us to consider a sequence of collections of information structures where the norm of covariance matrix is bounded below:

$$\mathcal{X}_n := \{X : \|\operatorname{Cov}(X)\|_2 \ge 1/n^2\},$$

and \mathcal{X}_n emulates the set of the posterior with the k-th sample where $k \leq n$. Additionally, as n goes to infinity \mathcal{X} contains all possible non constant information structure. Finally, the sequence is homogeneous and isotropic: whether an information structure X is in \mathcal{X}_n is independent of the position and direction so that translation and rotation do not affect the membership.

Formally, the average variance of the posterior with the n-th signal is $\mathbb{E}_{S(\leq n)}[(P(W|S^{(\leq n)} - P(W|S^{(\leq n-1)})^2],$ and equals $\mathrm{Var}[(P(W|S^{(\leq n)})] - \mathrm{Var}[(P(W|S^{(\leq n-1)})] = \Theta(1/n^2),$ because $\mathbb{E}_{S(\leq n-1)}[\mathbb{E}_{S(n)}[P(W|S^{(\leq n)})P(W|S^{(\leq n-1)})] = \mathbb{E}_{S(\leq n-1)}[P(W|S^{(\leq n-1)})^2].$

Our main results in this section show that a quadratic scoring rule

$$Q(x) := \frac{d}{d-1} ||x - c||^2 \text{ where } c := (1/d, \dots, 1/d)$$
(13)

is optimal for the sequence \mathcal{X}_n . The first result (Theorem 5.1) shows that the information gain Eq. (2) is of order $1/n^2$ for any convex function, and the constant is maximized if and only if H is the quadratic function Q. Then, we can translate the result into relative information gain, and show the quadratic scoring rule maximizes the relative information gain against all smooth function on $\operatorname{int}(\Delta_d)$ for $(\mathcal{X}_n)_{n\in\mathbb{Z}_{>0}}$ in Corollary 5.2.

Theorem 5.1. Let $\Omega = [d]$ and $H \in \mathcal{H} \cap \mathcal{C}^{\infty}(\operatorname{int}(\Delta_d))$.

$$\lim_{n \to \infty} n^2 \operatorname{Obj}_{\mathcal{X}_n}(H) \le \frac{d}{d-1},$$

and the equality holds if and only if H(x) = Q(x) for all $x \in \Delta_d$.

Corollary 5.2. For any finite set Ω with $|\Omega| = d$, and any smooth $H^* \in \mathcal{H}$, H^* maximizes the following relative information gain

$$R_{\mathrm{uni}}(H) := \inf_{\tilde{H} \in \mathcal{H} \cap \mathcal{C}^{\infty}(\mathrm{int}(\Delta_d))} \lim_{n \to \infty} \frac{\mathrm{Obj}_{\mathcal{X}_n}(H)}{\mathrm{Obj}_{\mathcal{X}_n}(\tilde{H})}.$$

if and only if $H^* = Q(x)$ for all $x \in \Delta$.

To prove Theorem 5.1, the main idea is that the information gain $\mathfrak{J}_H(X)$ only depends on the "curvature" of H at $\mathbb{E}X$ that is the strongly-convexity of H. (Lemma 5.4) Then, it is sufficient to show that the quadratic function Q has the largest possible curvature in \mathcal{H} . (Lemma 5.5)

The following lemma show it is sufficient to consider random variable X with $\|\operatorname{Cov}(X)\|_2 = 1/n^2$.

Lemma 5.3 (scaling). Let $H \in \mathcal{H}$ and $n \in \mathbb{Z}_{>0}$. For all $X \in \mathcal{X}_n$, there exists $\tilde{X} \in \mathcal{X}_n$ so that $\|\operatorname{Cov}(\tilde{X})\|_2 = 1/n^2$ and $\mathfrak{J}_H(X) \geq \mathfrak{J}_H(\tilde{X})$.

Then we connect information gain of \mathcal{X}_n to strongly convexity of a function. The proofs of Lemmas 5.3 and 5.4 are in Appendix C.

Lemma 5.4 (strongly convex). If $H \in \mathcal{H} \cap \mathcal{C}^{\infty}(\operatorname{int}(\Delta_d))$, H is α -strongly convex in $\operatorname{int}(\Delta)$ if and only if $\operatorname{inf}_{X \in \mathcal{X}_n: \|\operatorname{Cov}(X)\| = 1/n^2} \mathfrak{J}_H(X) \geq \frac{\alpha}{2n^2} + o\left(\frac{1}{n^2}\right)$.

Finally, we show the most "strongly convex" function in \mathcal{H} is the quadratic scoring rule Q. Intuitively, if a function H is more strongly convex than Q, the difference H - Q is a convex function, and H is greater than Q at one of the vertex \hat{e}_k which is a contradiction because $Q(\hat{e}_k) = 1$ for all $k \in [d]$.

Lemma 5.5. Let $H \in \mathcal{H} \cap \mathcal{C}^{\infty}(\operatorname{int}(\Delta_d))$. H is $\frac{2d}{d-1}$ -strongly convex on Δ if and only if for all $x \in \Delta$, H(x) = Q(x).

Proof of Lemma 5.5. If $H \in \mathcal{H}$ is $\frac{2d}{d-1}$ -strongly convex, we define the symmetrization of H as

$$H_{sym}: (x_1, \dots, x_d) \mapsto \frac{1}{d!} \sum_{\sigma \in S_d} H(x_{\sigma(1)}, \dots, x_{\sigma(d)})$$

where S_d is the set of all permutations on [d]. We first prove

$$H(c) = 0 \text{ and } \nabla H(c) = 0. \tag{14}$$

Since the symmetrization H_{sym} is also $\frac{2d}{d-1}$ -strongly convex, smooth, and in \mathcal{H} , the following function is convex on Δ , $K(x) := H_{sym}(x) - H_{sym}(c) - Q(x)$. Because H_{sym} is convex and symmetric with respect to c, $H_{sym}(c) = \min_{x \in \Delta} H_{sym}(x)$, and $\nabla H_{sym}(c) = 0$ which proves the second part of Eq. (14). For the first part of Eq. (14), $\nabla K(c) = \nabla H_{sym}(c) - \nabla Q(c) = 0$, and $K(x) \geq K(c) = 0$ for all $x \in \Delta$ by Eq. (5). As a result, for any vertex \hat{e}_j ,

$$1 \ge H_{sym}(\hat{e}_j) \ge H_{sym}(c) + Q(\hat{e}_j) \ge 1,$$

so $H_{sym}(c) = 0$, and the vertices $K(\hat{e}_j) = 0$ for all j = 1, ..., d. Because $K \ge 0$ is convex and the evaluations at vertices are zero, K(x) = 0 for all $x \in \Delta$, so

$$H_{sym}(x) = Q(x), \text{ for all } x \in \Delta.$$
 (15)

Finally, because H is smooth and $\frac{2d}{d-1}$ -strongly convex, by Lemma 3.3, for all $x \in \text{int}(\Delta)$ and $v \in \mathbb{R}^d$ with $\mathbf{1}^\top v = 0$, $v^\top \nabla^2 H(x) v \geq \frac{2d}{d-1} ||v||^2$. Thus

$$v^{\top} \nabla^2 H_{sym}(x) v = \frac{1}{d!} \sum_{\sigma} v^{\top} \nabla^2 H(x_{\sigma(1)}, \dots, x_{\sigma(d)}) v \ge \frac{1}{d!} \sum_{\sigma} \frac{2d}{d-1} ||v||^2 = \frac{2d}{d-1} ||v||^2.$$

However, by Eq. (15), $v^{\top}\nabla^2 H_{sym}(x)v \leq \frac{2d}{d-1}\|v\|^2$, and $v^{\top}\nabla^2 H(x)v = \frac{2d}{d-1}\|v\|^2$ for all $x \in \text{int}(\Delta)$. Thus, we have a second order differential equation with initial conditions in Eq. (14), and we get H(x) = Q(x) for all $x \in \Delta$.

Proof of Theorem 5.1. First Q is $\frac{2d}{d-1}$ -strongly convex and $Q \in \mathcal{H}$. By Lemmas 5.3 and 5.4, $\operatorname{Obj}_{\mathcal{X}_n}(Q) = \frac{d}{d-1}\frac{1}{n^2} + o(\frac{1}{n^2})$, so $\lim_{n \to \infty} n^2 \operatorname{Obj}_{\mathcal{X}_n}(Q) = \frac{d}{d-1}$.

Conversely, suppose $H \in \mathcal{H}$ is smooth and $\lim_{n\to\infty} n^2 \operatorname{Obj}_{\mathcal{X}_n}(H) = \frac{d}{d-1}$. Then $\operatorname{Obj}_{\mathcal{X}_n}(H) = \frac{d}{d-1} \frac{1}{n^2} + o(\frac{1}{n^2})$. By Lemma 5.4, H is $\frac{2d}{d-1}$ -strongly convex in $\operatorname{int}(\Delta)$. By Lemma 5.5, H = Q in Δ . This completes the proof.

5.2 Exponential Family Information Structures

Now we consider two common families of information structures in Bayesian inference: Beta-Bernoulli and Dirichlet-categorical information structures. Specifically, the principal wants to predict the outcome of a d-sided dice. First, the principal announces her proper scoring rule to encourage the agent to collect signals through rolling the dice once (first setting) or a fixed amount of times (second setting). The agent then reports his prediction. Finally, the dice is rolled and the agent gets his payment according to the outcome and the proper scoring rule.

In Section 5.2.1 we consider d=2. As an warm-up, we begin with the static settings where the agent can only collect a fixed number of samples and apply Theorem 4.4 to solve the optimal scoring rules. Then we study the asymptotic setting, and show the log scoring rule is optimal. We further extend this result to general d-sided dice in Section 5.2.2.

5.2.1 Beta-Bernoulli Information Structure

Suppose we want to collect predictions on an outcome of a coin. Given $p \in [0, 1]$, the outcome W follows the Bernoulli distribution Bern(p) with $\Pr[W = 1] = p$, whereas the true value of p is unknown. We consider the following two cases.

- 1. Knowing that the agent privately observes N i.i.d samples from the coin, we want to incentivize the agent to make one additional observation.
- 2. We want to incentivize the agent to collect N+1 samples adaptively. That is for any sequence of $n \leq N$ samples from the coin, the information gain of the (n+1)-th sample is large enough to incentivize the agent to observe.

In both cases, the agent starts with an uninformative prior on the parameter of the coin p that is a uniform distribution on [0,1].

To formalize this, let $\theta \in \Theta = [0,1]$, $\Omega = \mathcal{S} = \{0,1\}$, and $m \in \mathbb{N}_{>0}$. Let $Beta(\theta,m)$ be a Beta distribution where the probability density at $p \in [0,1]$ is proportional to $p^{\theta m-1}(1-p)^{(1-\theta)m-1}$. Here θ is the mean and m is the effective sample size. We define Beta-Bernoulli information structure with (θ,m) as follow: p is sampled from $Beta(\theta,m)$. W and S are sampled independently and identically from Bern(p).

Static Setting In the first case, if the agent already observes n_1 heads and $n-n_1$ tails, his prior on W=1 is $\theta = \frac{n_1+1}{n+2}$, and the posterior predictive distribution is $X_{\theta}^{(n)} := P(W=1 \mid S)$ that is a Beta-Bernoulli information structure with $(\theta, n+2)$,

$$X_{\theta}^{(n)} = \begin{cases} \frac{(n+2)\theta+1}{n+3} & \text{with probability } \theta\\ \frac{(n+2)\theta}{n+3} & \text{with probability } 1-\theta \end{cases}$$
 (16)

Note that because the agent starts with the uniform prior, his effective sample size after n observation is m = n + 2.

By Eq. (4), the information gain of the (n+1)-th observation under a convex function $H:[0,1]\to\mathbb{R}$ is $\mathfrak{J}_H(X_{\theta}^{(n)})=\mathbb{E}[H(X_{\theta}^{(n)})]-H(\theta)$. Thus to elicit the (N+1)-th observation, the objective value in Eq. (2) is

$$Obj_{\left\{X_{\theta}^{(N)}:\theta=\frac{1}{N+2},\dots,\frac{N+1}{N+2}\right\}}(H) = \inf_{\theta=\frac{1}{N+2},\dots,\frac{N+1}{N+2}} \mathfrak{J}_{H}(X_{\theta}^{(N)}). \tag{17}$$

When $N \in \mathbb{Z}_{\geq 0}$ is finite, the collection of information structures $\left\{X_{\theta}^{(N)}: \theta = \frac{1}{N+2}, \ldots, \frac{N+1}{N+2}\right\}$ is finite. Hence, finding an optimal H for Eq. (17) is studied in Section 4.2. For the second case, to incentivize the agent to collect N+1 costly observation adaptively, we need to ensure for any sequence of $n \leq N$ samples, the agent's expected information gain of the (n+1)-th observation is large. Thus, the collection of information structures is $\left\{X_{\theta}^{(n)}: n \leq N, \theta = \frac{1}{n+2}, \ldots, \frac{n+1}{n+2}\right\}$ which is also finite. Figure 2 shows examples of the optimal scoring rules in the above two cases with a finite effective sample size N.

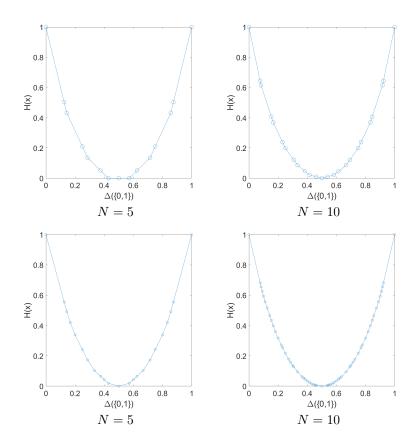


Figure 2: We use Theorem 4.4 to solve the optimal H for Beta-Bernoulli information structures with various N, and the top row is the first case and the bottom row is the second case.

However, as N increases the size of corresponding linear programming becomes larges. For instance, if we want to incentivize the agent to collect N+1 samples adaptively, the number of variable is $\Theta(N^2)$ which

is the length of Farev sequence of order N+1.⁷ In the next section, we are going to show the optimal scoring rules is the log scoring rule in Eq. (23) when $N \to \infty$.

Asymptotic Setting Instead of the static setting with a fixed number of samples, now we consider the agent can indefinitely get samples, and define a sequence of collections of Beta-Bernoulli information structures as follows.

$$\bar{\mathcal{X}}_{Beta}^{\leq N} := \left\{ X_{\theta}^{(n)} : n \leq N, \delta_N < \theta < 1 - \delta_N \right\}$$
(18)

which is the collection of all Beta-Bernoulli information structures with an effective sample size at most N and $\delta_N < \theta < 1 - \delta_N$ where $\delta_N \in o(1)$ and $\delta_N \in \omega(1/N)$. Note that $\bar{\mathcal{X}}_{Beta}^{\leq N}$ is an increasing sequence and $\lim_{N \to \infty} \bar{\mathcal{X}}_{Beta}^{\leq N} = \mathcal{X}_{Beta} := \{X_{\theta}^{(N)} : 0 < \theta < 1, N \in \mathbb{N}\}$ contains all Beta-Bernoulli information structures. Our main result in this section shows that a log scoring rule

$$H_{\ln}(x) := \frac{x \ln(x)}{\ln 2} + \frac{(1-x)\ln(1-x)}{\ln 2} + 1 \tag{19}$$

is optimal for the sequence $\bar{\mathcal{X}}_{Beta}^{\leq n}$

Theorem 5.6. Let $H \in \mathcal{H} \cap \mathcal{C}^{\infty}((0,1))$ be smooth on (0,1).

$$\lim_{N \to \infty} (N+3)^2 \operatorname{Obj}_{\bar{\mathcal{X}}_{Beta}^{\leq N}}(H) \leq \frac{1}{2 \ln 2}.$$

Moreover the equality holds if and only if $H(x) = H_{ln}(x)$ for all x in [0, 1].

We can also rewrite this result in terms of relative information gain, and show the log scoring rule maximizes the following relative information gain on the sequence in Eq. (18) against any function in $\mathcal{H} \cap$ $\mathcal{C}^{\infty}((0,1)).$

Corollary 5.7. For all $H^* \in \mathcal{H} \cap \mathcal{C}^{\infty}((0,1))$ smooth on (0,1), H^* maximizes the following relative information gain

$$R_{\mathrm{Beta}}(H) := \inf_{\tilde{H} \in \mathcal{H} \cap \mathcal{C}^{\infty}((0,1))} \lim_{n \to \infty} \frac{\mathrm{Obj}_{\tilde{\mathcal{X}}^{\leq n}_{Beta}}(H)}{\mathrm{Obj}_{\tilde{\mathcal{X}}^{\leq n}_{Beta}}(\tilde{H})},$$

if and only if $H^* = H_{ln}(x)$ for all $x \in (0,1)$.

In contrast to vanishing covariance information structures in Section 5.1, the variance of a Beta-Bernoulli information structure depends on the mean. For example, the posterior $X_{0.5}^{(N)}$ with mean 0.5 changes rapidly, but $X_{0.98}^{(N)}$ with mean 0.98 hardly changes. Additionally, the movement becomes small as the prior mean approaches the vertices 0 or 1, so a good scoring rule should be more curved as well.

The proof structure is similar to Theorem 5.1. We defer the proofs of the following two lemmas to Appendix D.

Lemma 5.8 (scaling). Let $H \in \mathcal{H}$. For all $\theta \in [0,1]$ and $n \leq N$, we have $\mathfrak{J}_H(X_{\theta}^{(n)}) \geq \mathfrak{J}_H(X_{\theta}^{(N)})$.

The next lemma relates the objective value to a second order differential operator $\mathfrak{D}_{\text{beta}}$

$$h(x) \mapsto \mathfrak{D}_{\text{beta}} h(x) = x(1-x)h''(x)$$

for any smooth function $h:(0,1)\to\mathbb{R}$ where h'' is the second derivative of h. Note that for the log scoring rule, $\mathfrak{D}_{\mathrm{beta}} H_{\ln}(x) = \frac{1}{\ln 2}$ for all $x \in (0, 1)$.

Lemma 5.9. For any smooth convex $H \in \mathcal{H} \cap \mathcal{C}^{\infty}((0,1))$, and $X \in \bar{\mathcal{X}}_{Beta}^N$,

$$\mathfrak{J}_H(X) = \frac{1}{2(N+3)^2} \, \mathfrak{D}_{\mathrm{beta}} \, H(\mathbb{E} \, X) + O\left(\frac{1}{(N+3)^3}\right)$$

when N is large enough. Moreover, $\text{Obj}_{\bar{\mathcal{X}}_{Beta}^N}(H_{\ln}) = \frac{1}{2(\ln 2)(N+3)^2} + o\left(\frac{1}{(N+3)^2}\right)$, for log scoring rule.

⁷ Farey sequence of order m is the sequence of completely reduced fractions between 0 and 1 and the denominators are less than or equal to m.

Informally, the first part of Lemma 5.9 is a "pointwise" characterization where for each θ , $X_{\theta}^{(N)}$ the information gain pointwise converges to $\frac{\mathfrak{D}_{\text{beta}}H(\theta)}{2(N+3)^2}$, and the moreover part is an "almost uniform" characterization, where the information gain converges to the limit almost uniformly in (0,1), because $(\delta_N, 1 - \delta_N) \to (0,1)$ as $N \to \infty$.

Finally, we show the log scoring rule with H_{ln} has the largest $\inf_x \mathfrak{D}_{beta} H(x)$ for all smooth $H \in \mathcal{H}$.

Lemma 5.10 (optimal). If $H \in \mathcal{H} \cap \mathcal{C}^{\infty}((0,1))$,

$$\inf_{x \in (0,1)} \mathfrak{D}_{\text{beta}} H(x) \le \frac{1}{\ln 2}.$$

Moreover, if $\inf_{x \in (0,1)} \mathfrak{D}_{\text{beta}} H(x) = \frac{1}{\ln 2}$, H equal the log scoring rule H_{ln} on (0,1).

Proof of Lemma 5.10. Because for all $x \in (0,1)$ $H''_{\ln}(x) = \frac{1}{(\ln 2)x(1-x)}$, $\mathfrak{D}_{\text{beta}} H_{\ln}(x) = 1/\ln 2$. Additionally $H_{\ln}(0) = H_{\ln}(1) = 1$.

Now we show the infimum cannot be greater than $1/\ln 2$ for any smooth $H \in \mathcal{H}$; otherwise $H = H_{ln}$. Suppose

$$\inf \mathfrak{D}_{\text{beta}} H \ge 1/\ln 2. \tag{20}$$

First by a similar argument in the proof of Lemma 5.5, with out lose of generality, we only need to consider H is symmetric about 0.5 and H(0.5) = 0. Let $K := H - H_{\rm ln}$ which is a smooth function on (0,1) and $K(0.5) = H(0.5) - H_{\rm ln}(0.5) = 0$. By Eq. (20), $\mathfrak{D}_{\rm beta} K(x) = x(1-x)K''(x) \geq 0$ for all $x \in (0,1)$ which yields $K''(x) \geq 0$. Because K is convex, the maximum happens at the boundary 0 or 1 and greater than K(0.5) = 0,

$$1 + K(0.5) \le 1 + \max\{K(0), K(1)\} = 1 + \max\{H(0) - H_{\ln}(0), H(1) - H_{\ln}(1)\} = \max\{H(0), H(1)\}.$$

Additionally, since $H \in \mathcal{H}$, $H(0), H(1) \leq 1$. Combining these two, we have K(0) = K(0.5) = K(1) = 0, so K(x) = 0 for all $x \in [0, 1]$. Therefore, $H(x) = H_{\ln}(x)$ for all x.

Now, let's put everything together.

Proof of Theorem 5.6. First, for any smooth $H \in \mathcal{H}$ not equal to H_{ln} , $\lim_{n\to\infty} 2(n+3)^2 \, \mathfrak{J}_H(X_{\theta}^{(n)}) = \mathfrak{D}_{\text{beta}} H(\theta)$ for all $\theta \in (0,1)$ by Lemma 5.9. By Lemma 5.10 and 5.8, we have

$$\lim_{N\to\infty} (N+3)^2 \operatorname{Obj}_{\bar{\mathcal{X}}_{beta}^{\leq N}}(H) \leq \inf_{\theta} \frac{1}{2} \mathfrak{D}_{beta} H(\theta) < \frac{1}{2 \ln 2}.$$

The converse holds by Lemma 5.9 and 5.8.

Note that for any N the set $\bar{\mathcal{X}}_{Beta}^{\leq N}$ does not contain all Beta-Bernoulli information structures with effective sample size at most N, but it additionally requires the mean $\delta_N \leq \theta \leq 1 - \delta_N$ is bounded away from the vertices, 0 and 1. This restriction is due to the log scoring rule's singular behavior around the vertices. In particular, in Lemma 5.9 the information gain $\mathfrak{J}_{H_{\rm in}}(X_{\theta}^{(n)})$ only converges pointwise to $\mathfrak{D}_{\rm beta}\,H_{\rm ln}(\theta)$ and the constant in the error term increases as θ approaches the vertices. In other words, Theorem 5.6 ensures the log scoring rule is optimal in an almost uniform manner, but Theorem 5.1 shows the quadratic scoring rule is optimal in a uniform manner. In Fig. 3, we numerically show the optimal scoring rules under the first case with finite N converge to the log scoring rule.

5.2.2 Dirichlet-Categorical Information Structure

We will show the log scoring rule is still the optimal scoring rule for the collection of Dirichlet-categorical information structures.

If $|\Omega| = d$, let $\theta \in \Theta = \Delta(\Omega)$ and $m \in \mathbb{N}_{>0}$. Let $Dir(\theta, m)$ be the Dirichlet distribution where the probability density at $p \in \Delta(\Omega)$ is proportional to $\Pi_{k=1}^d p_k^{\theta_k m-1}$. Let θ be the mean and m be the effective sample size. We similarly define Dirichlet-categorical information structure with (θ, m) as follow: p is sampled from $Dir(\theta, m)$. The state of the world W and signal S are sampled independently from Cat(p) where the outcome is k with probability p_k .

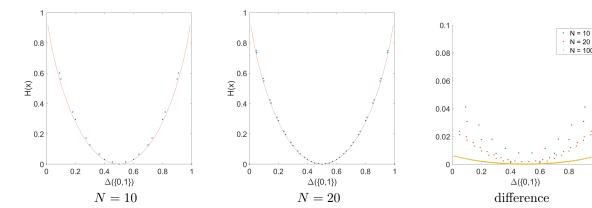


Figure 3: We plot the log scoring rules H_{ln} and the optimal scoring rules under the static setting which is the first case with constant N. The rightmost panel shows the difference between those optimal scoring rules and the log scoring rule.

We will generalize Beta-Bernoulli information structures (16) and the log scoring rule Eq. (19) with d=2to general d-dimensional cases. When the agent's prior is $Dir(\theta, n+d)$, the information structure of the (n+1)-th sample is $X_{\theta}^{(n)}$ where for all $k \in [d]$

$$\Pr\left[X_{\theta}^{(n)} = \frac{(n+d)}{n+d+1}\theta + \frac{1}{n+d+1}\hat{e}_k\right] = \theta_k.$$
 (21)

Furthermore, we consider the following sequence of collections of information structures

$$\bar{\mathcal{X}}_{dir}^{\leq N} := \left\{ X_{\theta}^{(n)} : n \leq N, \delta_N < \theta_k, \forall k \in [d] \right\}$$
 (22)

where $\delta_N \in \omega(1/N)$ and $\delta_N \in o(1)$.

Our main results (Theorem 5.11 and Corollary 5.12) show that a log scoring rule

$$H_{\ln}(x) = \frac{1}{\ln d} \sum_{k=1}^{d} x_k \ln x_k + 1,$$
(23)

is optimal for the sequence $\bar{\mathcal{X}}_{dir}^{\leq n}$. However, we only show H_{ln} is optimal against any function that is smooth on the close set Δ_d , instead of the relative interior $\operatorname{int}(\Delta_d)$. Notice that H_{ln} is smooth on $\operatorname{int}(\Delta_d)$ but not smooth on Δ_d , but several standard proper scoring rules are smooth on Δ_d , e.g., quadratic, and spherical scoring rule. Therefore, our results still provide relevant comparison among common scoring rules.

Theorem 5.11. Given $\Omega = [d]$, let $H \in \mathcal{H} \cap \mathcal{C}^{\infty}(\Delta_d)$ be smooth on Δ_d .

$$\lim_{N\to\infty}(N+d+1)^2\operatorname{Obj}_{\bar{\mathcal{X}}_{dir}^{\leq N}}(H)<\frac{d-1}{2\ln d}=\lim_{N\to\infty}(N+d+1)^2\operatorname{Obj}_{\bar{\mathcal{X}}_{dir}^{\leq N}}(H_{\ln}).$$

We can also rewrite this result in terms of relative information gain as Corollary 5.12.

Corollary 5.12. For all $H^* \in \mathcal{H} \cap \mathcal{C}^{\infty}(\Delta_d) \cup \{H_{ln}\}$, H^* maximizes the following relative information gain

$$R_{\mathrm{dir}}(H) := \sup_{\tilde{H} \in \mathcal{H} \cap \mathcal{C}^{\infty}(\Delta_d) \cup \{H_{\mathrm{ln}}\}} \lim_{n \to \infty} \frac{\mathrm{Obj}_{\bar{\mathcal{X}}_{dir}^{\leq n}}(\tilde{H})}{\mathrm{Obj}_{\bar{\mathcal{X}}_{dir}^{\leq n}}(H)},$$

if and only if $H^* = H_{\ln}(x)$ for all $x \in \operatorname{int}(\Delta_d)$.

The proof structure is mostly identical to Theorem 5.6. However, the main challenge is Lemma 5.15 which shows the log scoring rule maximizes a differential operator $\mathfrak{D}_{\rm dir}$ defined in Eq. (24). We prove Lemma 5.15 in Section 5.2.2, and omit rest of the proofs.

Lemma 5.13 (scaling). Let $H \in \mathcal{H}$. For all $\theta \in \Delta(\Omega)$ and $n \leq N$ we have $\mathfrak{J}_H(X_{\theta}^{(n)}) \geq \mathfrak{J}_H(X_{\theta}^{(N)})$.

The next lemma relates the information gain to a second order differential operator \mathfrak{D}_{dir} , for all smooth $h \in \mathcal{H}$ and $x \in \Delta(\Omega)$

$$h(x) \mapsto \mathfrak{D}_{dir} h(x) = \sum_{k=1}^{d} x_k (\hat{e}_k - x)^{\top} \nabla^2 h(x) (\hat{e}_k - x)$$
 (24)

where $(\hat{e}_k)_{k=1}^d$ are the vertices of simplex $\Delta(\Omega)$. The proof is basic identical to the proof of Lemma 5.9.

Lemma 5.14. For any smooth convex $H \in \mathcal{H} \cap \mathcal{C}^{\infty}(\operatorname{int}(\Delta_d))$, and Dirichlet-Categorical information structure $X = X_{\theta}^{(N)}$,

$$\mathfrak{J}_H(X) = \frac{1}{2(N+d+1)^2} \, \mathfrak{D}_{\mathrm{dir}} \, H(\mathbb{E} \, X) + O\left(\frac{1}{(N+d+1)^3}\right)$$

when N is large enough. Moreover, for log scoring rule, $\text{Obj}_{\bar{\mathcal{X}}_{Beta}^N}(H_{\ln}) = \frac{d-1}{2(\ln d)(N+d+1)^2} + o\left(\frac{1}{(N+d+1)^2}\right)$.

Finally, we show the infimum of \mathfrak{D}_{dir} on H_{ln} is greater than the infimum of \mathfrak{D}_{dir} on any $H \in \mathcal{H}$ which is smooth on Δ_d . The idea stems from the maximum principle of elliptic second order differential equations.

Lemma 5.15 (optimal). If $H \in \mathcal{H} \cap \mathcal{C}^{\infty}(\Delta_d)$,

$$\inf_{x \in \operatorname{int}(\Delta_d)} \mathfrak{D}_{\operatorname{dir}} H(x) < \frac{d-1}{\ln d}.$$

Moreover, $\inf_{x \in \operatorname{int}(\Delta_d)} \mathfrak{D}_{\operatorname{dir}} H_{\operatorname{ln}}(x) = \frac{d-1}{\ln d}$.

Proof of Lemma 5.15 Despite being similar to Lemma 5.5 and 5.10, the proof of Lemma 5.15 is quite challenging. To see this, Lemma 5.10 only needs to handle functions in an one dimensional space and a lower bound of $\mathfrak{D}_{\text{beta}}$ is sufficient to characterize such function fully. In particular, once we know $\mathfrak{D}_{\text{beta}} h(x)$ we can compute the second derivative of h at x. Similarly, Lemma 5.5 handles on functions on d-dimensional space, but the notion of strongly convex is also rich enough to fully determine the function. However, $\mathfrak{D}_{\text{dir}}$ is a differential operator with scalar value, so we cannot reconstruct h from the value of $\mathfrak{D}_{\text{dir}} h$. To address this, we use the idea of maximum principle (Theorem 3.5) to proof Lemma 5.15.

Note that $\mathfrak{D}_{\mathrm{dir}}(H_{\mathrm{ln}})(x) = \frac{d-1}{\ln d}$ for all x and $H_{\mathrm{ln}} \in [0,1]$. Suppose there exist a function $H \neq H_{\mathrm{ln}}$ in $\mathcal{H} \cap \mathcal{C}^{\infty}(\Delta_d)$ such that $\mathfrak{D}_{\mathrm{dir}}H(x) \geq \frac{d-1}{\ln d}$ for all $x \in \mathrm{int}(\Delta_d)$. We will show there exists an extreme point of the simplex \hat{e}_j such that

$$H(\hat{e}_i) > H_{ln}(\hat{e}_i) = 1,$$
 (25)

and prove Lemma 5.15 by contradiction.

By an argument similar to the proof of Lemma 5.5 and 5.10, we can assume that H is symmetric with respect to $c = (1/d, \ldots, 1/d)$, and H(c) = 0. Let $u := H - H_{ln}$. We have $u(c) = H(c) - H_{ln}(c) = 0$ and

$$\mathfrak{D}_{\operatorname{dir}} u(x) = \mathfrak{D}_{\operatorname{dir}} H(x) - \mathfrak{D}_{\operatorname{dir}} H_{\ln}(x) = \mathfrak{D}_{\operatorname{dir}} H(x) - \frac{d-1}{\ln d} \ge 0, \tag{26}$$

for all $x \in \text{int}(\Delta_d)$. Let $\theta^* \in \Delta_d$ be a maximum point of u. Now, we will show $\max u(x) > 0$ and $\hat{e}_1 = \theta^*$ is a maximum points which proves Eq. (25) and completes the proof.

Because \mathfrak{D}_{dir} in Eq. (24) is a second order differential operator with

$$A_{dir}(x) = \sum_{k} x_{k} (\hat{e}_{k} - x) (\hat{e}_{k} - x)^{\top}$$
(27)

where $A_{dir}(x)_{ii} = x_i(1-x_i)$ and $A_{dir}(x)_{ij} = -x_ix_j$, Eq. (26) is a partial differential equation on u. Ideally, we want to use the maximum principle (Theorem 3.5) and show that the maximum happens at the vertex so that $u(\hat{e}_k) > u(c) = 0$ which implies Eq. (25) and completes the proof. However, since $A_{dir}(x)\mathbf{1} = \mathbf{0}$ for all x, \mathfrak{D}_{dir} is not elliptic, and we cannot use the standard maximum principle. Thus we use a maximum principle which holds under this weaker condition.

Lemma 5.16. Let $V \subset \mathbb{R}^d$ be a linear subspace, and $W \subset \mathbb{R}^d$ be a bounded set where $\{x - w : \forall x \in W\} \subset V$ is relative open to V for some $w \in \mathbb{R}^d$. If a second order differential operator \mathcal{L} with A satisfies the following two conditions: 1) \mathcal{L} is uniformly elliptic on W which has $\rho > 0$ so that for all $x \in W$ and $v \in V$, $v^{\top}A(x)v \geq \rho ||v||_2^2$, 2) the image of A is in V, $null(A(x)) = V^{\perp}$ for all $x \in W$, then for any super-solution A satisfies A be A and continuous on A continu

$$\max_{cl(W)} h = \max_{bd(W)} h.$$

Moreover, if h has an interior maximum, h is a constant.

This result can be derived from Theorem 3.5, because we can use d-1 coordinate to encode the space. The first condition says \mathcal{L} is uniformly elliptic on the subspace and the second condition ensures the differential operator is well-defined. Alternatively this is also a special case of maximum principle on Riemannian manifolds. [26]. We include a proof in Appendix D for completeness.

Back to our proof of Lemma 5.15, though $\mathfrak{D}_{\text{dir}}$ is not uniformly elliptic on $\text{int}(\Delta_d)$, $\mathfrak{D}_{\text{dir}}$ is uniformly elliptic on $U_{\epsilon} := \{x \in \Delta_d : x_k \geq \epsilon, \forall k \in [d]\}$ for any $\epsilon > 0$. We set $\mathcal{V}_{dir} = \{v \in \mathbb{R}^d : \mathbf{1}^\top v = 0\}$ and $\{x - c : \forall x \in U_{\epsilon}\} \subset \mathcal{V}_{dir}$, and by Lemma 5.16 the maximum value of u happens at the boundary of U_{ϵ} . Suppose that u has an interior maximum θ^* , we can set ϵ small enough such that $\theta^* \in \text{int}(U_{\epsilon})$ which reaches a contradiction. Therefore, θ^* can only be in the boundary, and, due to the symmetry of u, we can assume

$$\theta^* \in \Delta_{d-1} \subset bd(\Delta(\Omega)) = bd(\Delta_d).$$

Since $u(c) = H(c) - H_{ln}(c) = 0$ and u is not constant, $u(\theta^*) > 0$.

Now we want to can the maximum principle recursively such that the maximum can be only at \hat{e}_1 which completes our proof. Formally, it is sufficient to prove $\sum_{k=0}^{d-1} x_k (\hat{e}_k - x)^\top \nabla^2 u(x) (\hat{e}_k - x) > 0$ for all $x \in \operatorname{int}(\Delta_{d-1})$ which is $\mathfrak{D}_{\operatorname{dir}}$ without the last coordinate. To prove this, $\lim_{x'\to x} \mathfrak{D}_{\operatorname{dir}} H(x') = \mathfrak{D}_{\operatorname{dir}} H(x) \geq \frac{d-1}{\ln d}$ since $H \in \mathcal{C}^{\infty}(\Delta(\Omega))$, and $\mathfrak{D}_{\operatorname{dir}} H_{\ln}(x) = \frac{d-2}{\ln d}$, by direct computation. Therefore, we can apply Lemma 5.16 again and show $\theta^* \in \Delta_{d-2}$. By induction, $\theta^* = \hat{e}_1$ and this completes our proof.

References

- [1] Jacob D. Abernethy, Yiling Chen, Chien-Ju Ho, and Bo Waggoner. Actively purchasing data for learning. CoRR, abs/1502.05774, 2015. URL http://arxiv.org/abs/1502.05774.
- [2] Shoshana Abramovich and Lars-Erik Persson. Some new estimates of the 'jensen gap'. *Journal of Inequalities and Applications*, 2016(1):1–9, 2016.
- [3] Anish Agarwal, Munther A. Dahleh, and Tuhin Sarkar. A marketplace for data: An algorithmic solution. CoRR, abs/1805.08125, 2018. URL http://arxiv.org/abs/1805.08125.
- [4] Luis J Alías, Paolo Mastrolia, and Marco Rigoli. *Maximum principles and geometric applications*, volume 700. Springer, 2016.
- [5] Yang Cai, Constantinos Daskalakis, and Christos H Papadimitriou. Optimum statistical estimation with strategic data sources. arXiv preprint arXiv:1408.2539, 2014.
- [6] Yiling Chen and Shuran Zheng. Prior-free data acquisition for accurate statistical estimation. CoRR, abs/1811.12655, 2018. URL http://arxiv.org/abs/1811.12655.
- [7] Yiling Chen, Nicole Immorlica, Brendan Lucier, Vasilis Syrgkanis, and Juba Ziani. Optimal data acquisition for statistical estimation. CoRR, abs/1711.01295, 2017. URL http://arxiv.org/abs/ 1711.01295.
- [8] Yiling Chen, Chara Podimata, Ariel D. Procaccia, and Nisarg Shah. Strategyproof linear regression in high dimensions, 2018.
- [9] Rachel Cummings, Stratis Ioannidis, and Katrina Ligett. Truthful linear regression. CoRR, abs/1506.03489, 2015. URL http://arxiv.org/abs/1506.03489.

- [10] Ofer Dekel, Felix Fischer, and Ariel D Procaccia. Incentive compatible regression learning. *Journal of Computer and System Sciences*, 76(8):759–777, 2010.
- [11] Lawrence C. Evans. Partial differential equations. American Mathematical Society, Providence, R.I., 2010. ISBN 9780821849743 0821849743.
- [12] Lisa K Fleischer and Yu-Han Lyu. Approximately optimal auctions for selling privacy when costs are correlated with data. In *Proceedings of the 13th ACM conference on electronic commerce*, pages 568–585, 2012.
- [13] Amirata Ghorbani and James Zou. Data shapley: Equitable valuation of data for machine learning, 2019.
- [14] Arpita Ghosh and Aaron Roth. Selling privacy at auction. In *Proceedings of the 12th ACM conference on Electronic commerce*, pages 199–208, 2011.
- [15] Arpita Ghosh, Katrina Ligett, Aaron Roth, and Grant Schoenebeck. Buying private data without verification. In *Proceedings of the fifteenth ACM conference on Economics and computation*, pages 931–948, 2014.
- [16] Malay Ghosh, Nitis Mukhopadhyay, and Pranab Kumar Sen. Sequential estimation, volume 904. John Wiley & Sons, 2011.
- [17] Richard D. Gill and Boris Y. Levit. Applications of the van Trees inequality: a Bayesian Cramér-Rao bound. Bernoulli, 1(1-2):59 79, 1995. doi: bj/1186078362. URL https://doi.org/.
- [18] Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation. Journal of the American Statistical Association, 102(477):359–378, 2007.
- [19] Moritz Hardt, Nimrod Megiddo, Christos Papadimitriou, and Mary Wootters. Strategic classification, 2015.
- [20] Jason D Hartline, Yingkai Li, Liren Shan, and Yifan Wu. Optimization of scoring rules. arXiv preprint arXiv:2007.02905, 2020.
- [21] Ronald A Howard. Information value theory. *IEEE Transactions on systems science and cybernetics*, 2 (1):22–26, 1966.
- [22] John McCarthy. Measures of the value of information. Proceedings of the National Academy of Sciences of the United States of America, 42(9):654, 1956.
- [23] Reshef Meir, Ariel D Procaccia, and Jeffrey S Rosenschein. Algorithms for strategyproof classification. Artificial Intelligence, 186:123–156, 2012.
- [24] Eric Neyman, Georgy Noarov, and S. Matthew Weinberg. Binary scoring rules that incentivize precision. CoRR, abs/2002.10669, 2020. URL https://arxiv.org/abs/2002.10669.
- [25] Kobbi Nissim, Salil P. Vadhan, and David Xiao. Redrawing the boundaries on purchasing data from privacy-sensitive individuals. *CoRR*, abs/1401.4092, 2014. URL http://arxiv.org/abs/1401.4092.
- [26] Pablo Padilla. The principal eigenvalue and maximum principle for second order elliptic operators on riemannian manifolds. *Journal of Mathematical Analysis and Applications*, 205(2):285–312, 1997.
- [27] Javier Perote and Juan Perote-Pena. Strategy-proof estimators for simple regression. *Mathematical Social Sciences*, 47(2):153–176, 2004.
- [28] Aaron Roth and Grant Schoenebeck. Conducting truthful surveys, cheaply. *CoRR*, abs/1203.0353, 2012. URL http://arxiv.org/abs/1203.0353.
- [29] Leonard J Savage. Elicitation of personal probabilities and expectations. *Journal of the American Statistical Association*, 66(336):783–801, 1971.

- [30] Bo Waggoner, Rafael Frongillo, and Jacob Abernethy. A market framework for eliciting private data. In *Proceedings of the 28th International Conference on Neural Information Processing Systems Volume 2*, NIPS'15, page 3510–3518, Cambridge, MA, USA, 2015. MIT Press.
- [31] Shuran Zheng, Bo Waggoner, Yang Liu, and Yiling Chen. Active information acquisition for linear optimization. *CoRR*, abs/1709.10061, 2017. URL http://arxiv.org/abs/1709.10061.

A Details and Proofs in Sections 2 and 3

Proof of Proposition 2.1. First note that the sequence $(\mathrm{Obj}_{X_n}(H))_n$ is non-increasing, because \mathcal{X}_1,\ldots is non-decreasing, and $\mathrm{Obj}_{X_n}(H)\geq 0$ for all n by the Jensen's inequality. Thus, $\lim_{n\to\infty}\mathrm{Obj}_{X_n}(H)$ exists. Now we show that

$$\lim_{n \to \infty} \mathrm{Obj}_{X_n}(H) = \mathrm{Obj}_{\mathcal{X}}(H) \text{ where } \mathcal{X} = \lim_{n \to \infty} X_n.$$
 (28)

Because $\mathcal{X}_n \subseteq \mathcal{X}$ for all n, $\lim_{n\to\infty} \mathrm{Obj}_{\mathcal{X}_n}(H) \ge \mathrm{Obj}_{\mathcal{X}}(H)$. To show $\lim_{n\to\infty} \mathrm{Obj}_{\mathcal{X}_n}(H) \le \mathrm{Obj}_{\mathcal{X}}(H)$, given any $\epsilon > 0$, there exists $X \in \mathcal{X}$ such that $\mathfrak{J}_H(X) \le \inf_{\tilde{X} \in \mathcal{X}} \mathfrak{J}_H(\tilde{X}) + \epsilon = \mathrm{Obj}_{\mathcal{X}}(H) + \epsilon$. Because $X \in \mathcal{X}$, there exists N so that for all $n \ge N$, $X \in \mathcal{X}_n$. Hence $\lim_n \mathrm{Obj}_{\mathcal{X}_n}(H) \le \mathfrak{J}_H(X)$. Combining these two, we have $\lim_n \mathrm{Obj}_{\mathcal{X}_n}(H) \le \mathrm{Obj}_{\mathcal{X}}(H) + \epsilon$ for all ϵ which proves Eq. (28). Therefore, if $\mathrm{Opt}(\mathcal{X}) > 0$, we can apply the limit operator separately in Eq. (3), and have

$$\inf_{\tilde{H}\in\mathcal{H}}\lim_{n\to\infty}\frac{\mathrm{Obj}_{\mathcal{X}_n}(H)}{\mathrm{Obj}_{\mathcal{X}_n}(\tilde{H})}=\inf_{\tilde{H}\in\mathcal{H}}\frac{\mathrm{Obj}_{\mathcal{X}}(H)}{\mathrm{Obj}_{\mathcal{X}}(\tilde{H})}=\frac{\mathrm{Obj}_{\mathcal{X}}(H)}{\mathrm{Opt}(\mathcal{X})}.$$

which shows the Problem 1 and 3 are equivalent.

Proof of Lemma 3.3. Suppose H is α -strongly convex in $\operatorname{int}(\Delta)$. Given any $x \in \operatorname{int}(\Delta)$ and $v \in \mathbb{R}^d$ with $\mathbf{1}^\top v = 0$, we can pick a small $t \in \mathbb{R}_{>0}$ such that $y := x + tv \in \operatorname{int}(\Delta)$. First by Taylor expansion, there exists some $0 \le \tau \le t$ such that

$$H(y) = H(x + tv) = H(x) + t\nabla H(x)^{\top} v + \frac{t^2}{2} v^{\top} \nabla^2 H(x + \tau v) v.$$

Then we apply Eq. (5), and have $v^{\top}\nabla^2 H(x+\tau v)v \geq \alpha ||v||^2$. Because H is smooth, by taking t to zero, we have

$$v^{\top} \nabla^2 H(x) v \ge \alpha ||v||^2$$
.

Conversely, if for all $x \in \operatorname{int}(\Delta)$ and $v \in \mathbb{R}^d$ with $\mathbf{1}^\top v = 0$, we have $v^\top \nabla^2 H(x) v \ge \alpha \|v\|^2$, we want to show Eq. (5) is true. For any $x, y \in \operatorname{int}(\Delta)$, we set $v = y - x \in \mathbb{R}^d$. Again by Taylor's expansion, there exists some $z \in \operatorname{int}(\Delta)$ such that $H(y) = H(x) + \nabla H(x)^\top v + \frac{1}{2} v^\top \nabla^2 H(z) v$. Finally, because $x, y \in \operatorname{int}(\Delta)$ and $z \in \operatorname{int}(\Delta)$, $\sum v_j = 0$ and $v^\top \nabla^2 H(z) v \ge \alpha \|v\|^2$. Therefore, we complete the proof.

B Details and Proofs in Section 4.3

B.1 Proof of Proposition 4.5

We prove it by contradiction. Given a bounded H, let \mathcal{X}_H settles H.

Suppose the minimum of H is not zero. Firstly, if $\min H = 1$, then H is a constant and can never be uniquely optimal. Secondly, if $0 < \min H < 1$, we can define a new bounded convex function $\tilde{H}(x) := \frac{1}{1-\min H}(H(x)-1)+1$ for all $x \in \Delta_d$. Moreover, for any $X \in \mathcal{X}_H$, $\mathfrak{J}_{\tilde{H}}(X) = \frac{1}{1-\min H}\mathfrak{J}_H(X) > \mathfrak{J}_H(X)$ because $\mathfrak{J}_H(X) \geq \operatorname{Opt}(\mathcal{X}_H) > 0$ and $1-\min H < 1$. Therefore, $\operatorname{Obj}_{\mathcal{X}_H}(\tilde{H}) \geq \operatorname{Obj}_{\mathcal{X}_H}(H)$ and $\tilde{H} \neq H$, so H cannot be uniquely optimal for \mathcal{X}_H which is a contradiction.

For the second condition, suppose there exists $j^* \in [d]$ such that $H(\hat{e}_{j^*}) < 1$. Because \hat{e}_j for $j \in [d]$ are all the vertices of Δ_d , any $x \in \Delta_d$ is a convex combination of those vertices, $x = \sum_j \alpha_j \hat{e}_j$ where $\alpha_j \geq 0$ and $\sum_j \alpha_j = 1$. With this notion, we define an affine function $K(x) = K(\sum_j \alpha_j \hat{e}_j) := \sum_j \alpha_j (1 - H(\hat{e}_j))$, and two new bounded convex functions $H_1(x) := H(x) + K(x)$ and $H_2(x) = \frac{1}{1 - \min H_1} (H_1(x) - 1) + 1$ for all $x \in \Delta_d$. First, because adding the affine function K does not change the information gain, $\operatorname{Obj}_{\mathcal{X}_H}(H_1) = \operatorname{Obj}_{\mathcal{X}_H}(H)$. Additionally, because $\min H_1 > 0$, $\operatorname{Obj}_{\mathcal{X}_H}(H_1) < \operatorname{Obj}_{\mathcal{X}_H}(H_2)$. Combining these two inequality, we have $\operatorname{Obj}_{\mathcal{X}_H}(H) < \operatorname{Obj}_{\mathcal{X}_H}(H_2)$.

B.2 Upper and Lower Bound Evaluations of a Convex Function

We prove this result by constructing a set of information structures to upper bound (Lemma B.1) and lower bound (Lemma B.2) the evaluation of an optimal convex function at a point $\theta \in \Delta_d$.

Lemma B.1 (upper bounds). Let H satisfy Eq. (12) and \mathcal{X} be a collection of information structures which settles H and $0 < \operatorname{Opt}(\mathcal{X})$. For any $\theta \in \Delta_d$, there exists a constant $C_\theta > 0$ so that if $\operatorname{Opt}(\mathcal{X}) \leq C_\theta$, there exists an information structure X_θ^{up} such that $\mathcal{X} \cup \{X_\theta^{up}\}$ settles H and $H^*(\theta) \leq H(\theta)$ for any optimal H^* for $\mathcal{X} \cup \{X_\theta^{up}\}$.

Proof of Lemma B.1. Because Δ_d is a convex hull of $(\hat{e}_j)_{j \in [d]}$, we use $\alpha_j(x)$ to denote the weight of vertex \hat{e}_j for the point x such that $x = \sum_j \alpha_j(x)\hat{e}_j$.

Given $\theta \in \Delta_d$, if $H(\theta) = 1$, we can pick an arbitrary information structure in \mathcal{X} , and the two conditions trivially holds. Otherwise, if $H(\theta) < 1$ and $\operatorname{Opt}(\mathcal{X}) \leq C_{\theta} := 1 - H(\theta)$, we construct X_{θ}^{up} such that $\Pr[X_{\theta} = \theta] = 1 - \frac{\operatorname{Opt}(\mathcal{X})}{1 - H(\theta)}$, and $\Pr[X_{\theta} = \hat{e}_j] = \frac{\operatorname{Opt}(\mathcal{X})}{1 - H(\theta)} \alpha_j(\theta)$ for all $j \in [d]$. First, let's compute the information gain of X_{θ}^{up} on H,

$$\mathfrak{J}_{X_{\theta}^{up}}(H) = \sum_{j} \frac{\operatorname{Opt}(\mathcal{X})}{1 - H(\theta)} \alpha_{j}(\theta) H(\hat{e}_{j}) + \left(1 - \frac{\operatorname{Opt}(\mathcal{X})}{(1 - H(\theta))}\right) H(\theta) - H(\mathbb{E} X_{\theta}^{up}) \\
= \frac{\operatorname{Opt}(\mathcal{X})}{1 - H(\theta)} + \left(1 - \frac{\operatorname{Opt}(\mathcal{X})}{1 - H(\theta)}\right) H(\theta) - H(\mathbb{E} X_{\theta}) \qquad (H(\hat{e}_{j}) = 1 \text{ for all } j) \\
= \frac{\operatorname{Opt}(\mathcal{X})}{1 - H(\theta)} - \frac{\operatorname{Opt}(\mathcal{X})}{1 - H(\theta)} H(\theta) = \operatorname{Opt}(\mathcal{X}) \qquad (\mathbb{E} X_{\theta} = \theta)$$

Thus, $\operatorname{Obj}_{\mathcal{X} \cup \{X_{\theta}^{up}\}}(H) = \operatorname{Opt}(\mathcal{X}) = \operatorname{Opt}(\mathcal{X} \cup \{X_{\theta}^{up}\})$. When H^{\star} is optimal for $\mathcal{X} \cup \{X_{\theta}^{up}\}$), $\mathfrak{J}_{X_{\theta}^{up}}(H^{\star}) \geq \operatorname{Opt}(\mathcal{X})$, and $1 - H^{\star}(\theta) \geq 1 - H(\theta)$ which proves the second condition.

If θ^* is a minimum point of H, by Lemma B.1, we set X^* be the information structure so that $H^*(\theta^*) = H(\theta^*) = 0$ when X^* settles H^* . Lemma B.2 adds X^* to lower bound an evaluation of an optimal function.

Lemma B.2 (lower bounds). Let H satisfy Eq. (12) with minimum at θ^* and \mathcal{X} be a collection of information structures which settles H and $0 < \operatorname{Opt}(\mathcal{X})$. For any $\theta \in \Delta_d$, there exists a constant $D_{\theta} > 0$ so that if $\operatorname{Opt}(\mathcal{X}) \leq D_{\theta}$, there exist information structures X_{θ}^{lo} and X^* such that $\mathcal{X} \cup \{X_{\theta}^{lo}, X^*\}$ settles H and $H^*(\theta) \geq H(\theta)$ for any H^* that is optimal for $\mathcal{X} \cup \{X_{\theta}^{lo}, X^*\}$.

Proof of Lemma B.2. Since the convex function H satisfying Eq. (12), there exists $\theta^* \in \Delta_d$ such that $H(\theta^*) = 0$. First, if $H(\theta) = 0$, we can pick an arbitrary information structure in \mathcal{X} , and the two conditions trivially holds. If $H(\theta) > 0$, $\theta \neq \theta^*$ and there exists $j^* \in \arg\max_j \theta_j - \theta_j^*$. Hence, we can represent the minimum point θ^* as a convex combination of θ and \hat{e}_j with $j \in [d] \setminus \{j^*\}$, $\theta^* = \beta_{j^*}\theta + \sum_{j \neq j^*} \beta_j \hat{e}_j$. We define an information structure X_{θ}^{lo} such that $\Pr[X_{\theta}^{lo} = \theta] = \frac{\operatorname{Opt}(\mathcal{X})}{1 - \beta_{j^*}(1 - H(\theta))}\beta_{j^*}$, $\Pr[X_{\theta}^{lo} = \hat{e}_j] = \frac{\operatorname{Opt}(\mathcal{X})}{1 - \beta_{j^*}(1 - H(\theta))}\beta_j$ where $j \in [d] \setminus \{j^*\}$, and $\Pr[X_{\theta}^{lo} = \theta^*] = 1 - \frac{\operatorname{Opt}(\mathcal{X})}{1 - \beta_{j^*}(1 - H(\theta))}$ which is a valid distribution when $\operatorname{Opt}(\mathcal{X}) \leq D_{\theta} := 1 - \beta_{j^*}(1 - H(\theta))$. For the first condition, by direct computation,

$$\begin{split} &\mathfrak{J}_{X_{\theta}^{lo}}(H) \\ &= \frac{\operatorname{Opt}(\mathcal{X})}{1 - \beta_{j^*}(1 - H(\theta))} \left(\beta_{j^*}H(\theta) + \sum_{j \neq j^*} \beta_j H(\hat{e}_j)\right) + \left(1 - \frac{\operatorname{Opt}(\mathcal{X})}{1 - \beta_{j^*}(1 - H(\theta))}\right) H(\theta^*) - H(\mathbb{E}X_{\theta}^{lo}) \\ &= \frac{\operatorname{Opt}(\mathcal{X})}{1 - \beta_{j^*}(1 - H(\theta))} \left(\beta_{j^*}H(\theta) + 1 - \beta_{j^*}\right) - H(\mathbb{E}X_{\theta}^{lo}) \\ &= \frac{\operatorname{Opt}(\mathcal{X})}{1 - \beta_{j^*}(1 - H(\theta))} \left(\beta_{j^*}H(\theta) + 1 - \beta_{j^*}\right) = \operatorname{Opt}(\mathcal{X}) \end{split}$$

$$(\mathbb{E}X_{\theta}^{lo} = \theta^*)$$

Thus, $\operatorname{Obj}_{\mathcal{X} \cup \{X_{\theta}^{lo}\}}(H) = \operatorname{Opt}(\mathcal{X}) = \operatorname{Opt}(\mathcal{X} \cup \{X_{\theta}^{lo}\})$. Finally, by Lemma B.1 we can add X^* without changing the optimal value, and ensure $H^*(\theta^*) = 0$ for any optimal H^* . When H^* is optimal for $\mathcal{X} \cup \{X_{\theta}^{lo}\}$, $\mathfrak{I}_{X_{\theta}^{lo}}(H^*) \geq \operatorname{Opt}(\mathcal{X})$, and $\beta_{j^*}H^*(\theta) + 1 - \beta_{j^*} \geq \beta_{j^*}H(\theta) + 1 - \beta_{j^*}$ which proves the second condition. \square

B.3 Proof of Theorem 4.6

Now we use these lemmas to prove Theorem 4.6. To illustrate the idea, we first consider H is a v-shaped convex function on [0,1] with the minimum point at θ^* . We pick two points θ_0 and θ_1 with $0 < \theta_0 < \theta^* < \theta_1 < 1$. First we add an information X^* to \mathcal{X}_H such that $\Pr[X^* = \theta^*] = 1 - \epsilon$, $\Pr[X^* = 0] = (1 - \theta^*)\epsilon$, and $\Pr[X^* = 1] = \theta^*\epsilon$. For any $\epsilon > 0$, we have

$$H^{\star}(0) = H^{\star}(1) = 1 \text{ and } H^{\star}(\theta^{*}) = 0$$
 (29)

if H^* is optimal for $\{X^*\}$. If we take ϵ small enough, we can apply Lemmas B.1 and B.2 for θ_0 and θ_1 , and add four additional finite information structures to $\{X^*\}$ such that

$$H^{\star}(\theta_0) = H(\theta_0) \text{ and } H^{\star}(\theta_1) = H(\theta_1).$$
 (30)

Now we want to show the resulting finite collection of information structures \mathcal{X}_H ensures $H^* = H$ when H^* is optimal for \mathcal{X}_H . For all $0 < \theta < \theta_0$, by Eq. (29), H^* is a convex function below H because the epigraph of H is polytope with vertices at 0,1 and θ^* , so $H^*(\theta) \leq H(\theta)$. On the other hand, because $\theta_0 = \lambda \theta + (1 - \lambda)\theta^*$ for some $\lambda \in (0, 1)$, by Jensen's inequality $\lambda H^*(\theta) \geq H^*(\theta_0) - (1 - \lambda)H^*(\theta^*)$. Then by Eq. (30), $\lambda H^*(\theta) \geq H(\theta_0) - (1 - \lambda)H(\theta^*)$. Finally since H is affine in $[\theta_0, \theta^*]$, $H(\theta_0) - (1 - \lambda)H(\theta^*) = \lambda H(\theta)$. Therefore, $H^*(\theta) = H(\theta)$ for all $0 < \theta < \theta_1$. We can repeat this argument to the other three intervals where H is affine (θ_1, θ^*) , (θ^*, θ_1) , and $(\theta_1, 1)$.

The above proof can be easily adapted to any convex piecewise linear function.

Proof of Theorem 4.6. We prove it by constructing a finite collection of information structures \mathcal{X}_H .

Before defining \mathcal{X}_H , we introduce some notations. Because H satisfies Eq. (12), $H(\hat{e}_j) = 1$ for all $j \in [d]$ and there exists $\theta^* \in \Delta_d$ with $H(\theta^*) = 0$. Since $H : \Delta_d \to \mathbb{R}$ is a convex piecewise linear function, there exists a finite set of affine functions $A : \Delta_d \to \mathbb{R}$ and $A \in \mathcal{A}$, such that

$$H(x) = \max_{A \in \mathcal{A}} \left\{ A(x) \right\}.$$

Because H is convex and piecewise linear, the epigraph is a polytope. We define $F_A = \{x \in \Delta : H(x) = A(x)\}$ which is the set of points x where H is agree with the affine function A. Because F_A is also the domain of a face of epigraph of H, F_A is a convex of a finite number of vertices $v_1^A, \ldots v_{d_A}^A$. Finally we additionally pick an interior point $u^A \in F_A$ for each $A \in \mathcal{A}$.

Now, we are ready to define \mathcal{X}_H which is similar to the v-shaped case. First, we add X^* such that $\Pr[X^* = \theta^*] = 1 - \epsilon$, and $\Pr[X^* = \hat{e}_j] = (1 - \epsilon)\alpha_j(\theta^*)$ for all $j \in [d]$. For any $\epsilon > 0$, we have

$$H^{\star}(\theta^*) = 0 \text{ and } H^{\star}(\hat{e}_j) = 1 \text{ for all } j \in [d]$$
 (31)

if H^* is optimal for $\{X^*\}$. Because the set of vertices $\{v_k^A: A \in \mathcal{A}, k \leq d_A\}$ and interior points $\{u^A: A \in \mathcal{A}\}$ are finite, we can pick $\epsilon > 0$ small enough, and apply Lemmas B.1 and B.2 adding a finite number of finite information structures to \mathcal{X}_H such that for all $A \in \mathcal{A}, 1 \leq k \leq d_A$,

$$H^{\star}(v_k^A) = H(v_k^A) \tag{32}$$

and

$$H^{\star}(u^A) = H(u^A) \tag{33}$$

if H^* is optimal for \mathcal{X}_H .

For any $A \in \mathcal{A}$ and $\theta \in F_A$, by Eqs. (31) and (32) $H^*(\theta) \leq H(\theta)$. By Eqs. (32) and (33), $H^*(\theta) \geq H(\theta)$ which completes the proof.

B.4 Proof of Proposition 4.7

If $\theta \in \Delta_d$ with $\delta < H(\theta) < 1 - \delta$, we add the following two information structures:

•
$$X_{\theta}^{up}$$
: $\Pr[X_{\theta}^{up} = \theta] = 1 - \epsilon_{\theta}^{up}$ and $\Pr[X_{\theta}^{up} = \hat{e}_j] = \epsilon_{\theta}^{up} \alpha_j(\theta)$ where $\epsilon_{\theta}^{up} := \frac{\delta}{1 - H(\theta)}$ and $\theta = \sum_j \alpha_j(\theta) \hat{e}_j$.

• X_{θ}^{lo} : $\Pr[X_{\theta}^{lo} = \theta^*] = 1 - \epsilon_{\theta}^{lo}$, $\Pr[X_{\theta}^{lo} = \hat{e}_j] = \epsilon_{\theta}^{lo}\beta_j$ for all $j \neq j^*$, and $\Pr[X_{\theta}^{lo} = \theta] = \beta_{j^*}$ where $\epsilon_{\theta}^{lo} := \frac{\delta}{1 - \beta_{j^*}(1 - H(\theta))}$ and $\theta^* = \sum_{j \neq j^*} \beta_j \hat{e}_j + \beta_{j^*} \theta$.

Finally, we add an information structure X^* such that $\Pr[X^* = \theta^*] = 1 - \delta$ and $\Pr[X^* = \hat{e}_j] = \delta \alpha_j(\theta^*)$. By a direct computation, we have

$$\mathfrak{J}_H(X_{\theta}^{up}) = \mathfrak{J}_H(X_{\theta}^{lo}) = \mathfrak{J}_H(X^*) = \delta$$

for all θ with $\delta < H(\theta) < 1 - \delta$. Suppose H^* is optimal for $\mathcal{X}_{H,\delta} := \{X^{up}_{\theta}, X^{lo}_{\theta} : \delta < H(\theta) < 1 - \delta\} \cup \{X^*\}$. By the definition of X^{up}_{θ} , $H^*(\theta) \le H(\theta)$, and by the definition of X^{lo}_{θ} , $H^*(\theta) \ge H(\theta)$. Therefore, $H^*(\theta) = H(\theta)$ all θ with $\delta < H(\theta) < 1 - \delta$.

C Details and Proofs in Section 5.1

Lemma 5.3 (scaling). Let $H \in \mathcal{H}$ and $n \in \mathbb{Z}_{>0}$. For all $X \in \mathcal{X}_n$, there exists $\tilde{X} \in \mathcal{X}_n$ so that $\|\operatorname{Cov}(\tilde{X})\|_2 = 1/n^2$ and $\mathfrak{J}_H(X) \geq \mathfrak{J}_H(\tilde{X})$.

Proof of Lemma 5.3. Given $X \in \mathcal{X}_n$, we can decompose $X = \mu + D$ where $\mu = \mathbb{E} X \in \operatorname{int}(\Delta)$. Let $\tilde{X} = \mu + tD$ where $t = (n\sqrt{\|\operatorname{Cov}(X)\|_2})^{-1} \le 1$. Note that $\|\operatorname{Cov}(\tilde{X})\|_2 = t^2 \|\operatorname{Cov}(X)\|_2 = 1/n^2$, so $\tilde{X} \in \mathcal{X}_n$. Now it is sufficient to show

$$\mathfrak{J}_H(X) \ge \frac{1}{t} \, \mathfrak{J}_H(\tilde{X}) \ge \mathfrak{J}_H(\tilde{X}).$$
 (34)

To show this let T be a random variable that is independent of D and T=1 with probability t and 0 otherwise. We define $\check{X}=\mu+TD$. First we have

$$\mathfrak{J}_H(\check{X}) = \mathbb{E}_{D,T}[H(\mu + TD)] - H(\mu) = t\left(\mathbb{E}_D[H(\mu + D)] - H(\mu)\right) = t\mathfrak{J}_H(X).$$

On the other hand, by Jensen's inequality

$$\mathfrak{J}_{H}(\check{X}) = \mathbb{E}_{D,T}[H(\mu + TD)] - H(\mu)$$

$$\geq \mathbb{E}_{D}[H(\mathbb{E}_{T}[\mu + TD])] - H(\mu)$$

$$= \mathbb{E}_{D}[H(\mu + tD)] - H(\mu)$$

$$= \mathfrak{J}_{H}(\check{X}).$$

Combining these two, we prove Eq. (34) and the lemma.

Lemma 5.4 (strongly convex). If $H \in \mathcal{H} \cap \mathcal{C}^{\infty}(\operatorname{int}(\Delta_d))$, H is α -strongly convex in $\operatorname{int}(\Delta)$ if and only if $\operatorname{inf}_{X \in \mathcal{X}_n : \|\operatorname{Cov}(X)\| = 1/n^2} \mathfrak{J}_H(X) \ge \frac{\alpha}{2n^2} + o\left(\frac{1}{n^2}\right)$.

Proof of Lemma 5.4. Suppose H is α -strongly convex. Let $\mathbb{E} X = \mu$, we have

$$\begin{split} &\mathfrak{J}_H(X) \\ &= \mathbb{E}\left[H(X)\right] - H(\mathbb{E}\,X) \\ &\geq \mathbb{E}\left[H(\mu) + \nabla H(\mu)^\top (X - \mu) + \frac{\alpha}{2}\|X - \mu\|^2\right] - H(\mu) \qquad \qquad (H \text{ is α-strongly convex}) \\ &= \frac{\alpha}{2}\operatorname{Tr}(\operatorname{Cov}(X)) \qquad \qquad (\|X - \mu\|^2 = \operatorname{Tr}((X - \mu)(X - \mu)^\top)) \\ &\geq \frac{\alpha}{2}\|\operatorname{Cov}(X)\|_2 \qquad \qquad (\operatorname{Cov}(X) \text{ is positive semidefinite}) \\ &\geq \frac{\alpha}{2n^2} \end{split}$$

Conversely, if H is not α -strongly convex, by Lemma 5.3, there exists $x^* \in \operatorname{int}(\Delta)$, $\delta > 0$ and $v \in \mathbb{R}^d$ such that ||v|| = 1, $\sum v_j = 0$, and

$$v^{\top} \nabla^2 H(x^*) v < (\alpha - \delta) \|v\|^2. \tag{35}$$

We set X^* be a random variable such that $X^* = x^* + \frac{1}{n}v$ with probability 1/2 and $x^* - \frac{1}{n}v$ otherwise. Now we show X^* is in \mathcal{X}_n and $\mathfrak{J}_H(X^*)$ is smaller than $\alpha/(2n^2)$. First, when n is large enough, $\operatorname{supp}(X^*) \subset \operatorname{int}(\Delta)$.

 $\mathbb{E} X^* = x^* \text{ and } \operatorname{Cov}(X^*) = \mathbb{E}[(X^* - x^*)(X^* - x^*)^\top] = n^{-2}vv^\top, \text{ so } \|\operatorname{Cov}(X^*)\|_2 = 1/n^2. \text{ Thus } X^* \in \mathcal{X}_n.$ Second

$$\begin{split} \mathfrak{J}_{H}(X^{*}) = & \frac{1}{2} \left(H(x^{*} + \frac{1}{n}v) + H(x^{*} - \frac{1}{n}v) \right) - H(x^{*}) \\ = & \frac{1}{2n^{2}} v^{\top} \nabla^{2} H(x^{*})v + O(1/n^{3}) \qquad \text{(Taylor's expansion and } \{x : \|x - x^{*}\| \leq \frac{1}{n}\} \text{ is compact)} \\ < & \frac{(\alpha - \delta)}{2n^{2}} + O(1/n^{3}) \qquad \text{(by Eq. (35))} \end{split}$$

This completes the proof.

D Details and Proofs in Section 5.2

Lemma 5.8 (scaling). Let $H \in \mathcal{H}$. For all $\theta \in [0,1]$ and $n \leq N$, we have $\mathfrak{J}_H(X_{\theta}^{(n)}) \geq \mathfrak{J}_H(X_{\theta}^{(N)})$.

Proof of Lemma 5.8. Given a Beta-Bernoulli information structure X with (θ, n) , if $X \in \mathcal{X}_{Beta}^{\leq N}$, we have $n \leq N$. We want to show another Beta-Bernoulli information structure $\tilde{X} \in \mathcal{X}_{Beta}^{N}$ with (θ, N) has a smaller information gain.

We will reuse the proof of Lemma 5.3. Indeed, by direct computation, we have $\tilde{X} = \mathbb{E} X + \frac{n+3}{N+3}(X - \mathbb{E} X)$. Then by Eq. (34) with t = (n+3)/(N+3),

$$\mathfrak{J}(X) \ge \frac{N+3}{n+3}\,\mathfrak{J}(\tilde{X}) \ge \mathfrak{J}(\tilde{X}).$$

Thus, we complete the proof.

Lemma 5.9. For any smooth convex $H \in \mathcal{H} \cap \mathcal{C}^{\infty}((0,1))$, and $X \in \bar{\mathcal{X}}_{Beta}^N$

$$\mathfrak{J}_H(X) = \frac{1}{2(N+3)^2} \, \mathfrak{D}_{\mathrm{beta}} \, H(\mathbb{E} \, X) + O\left(\frac{1}{(N+3)^3}\right)$$

when N is large enough. Moreover, $\text{Obj}_{\bar{\mathcal{X}}_{Beta}^N}(H_{\ln}) = \frac{1}{2(\ln 2)(N+3)^2} + o\left(\frac{1}{(N+3)^2}\right)$, for log scoring rule.

Proof of Lemma 5.9. Let $x = \mathbb{E} X$.

$$\begin{split} &\mathfrak{J}_{H}(X) \\ &= \mathbb{E}\left[H(X)\right] - H(\mathbb{E}\,X) \\ &= xH\left(\frac{(N+2)x+1}{N+3}\right) + (1-x)H\left(\frac{(N+2)x}{N+3}\right) - H(x) \\ &= x\left(H\left(\frac{(N+2)x+1}{N+3}\right) - H(x)\right) + (1-x)\left(H\left(\frac{(N+2)x}{N+3}\right) - H(x)\right) \end{split}$$

Now we can estimate the above two term separately. Since H is smooth by Taylor's expansion, there exists ξ and $\tilde{\xi}$ such that $x \leq \xi \leq \frac{(N+2)x+1}{N+3}$, $\frac{(N+2)x}{N+3} \leq \tilde{\xi} \leq x$, and

$$H\left(\frac{(N+2)x+1}{N+3}\right) - H(x) = H'(x) \cdot \frac{1-x}{N+3} + \frac{1}{2}H''(x) \cdot \frac{(1-x)^2}{(N+3)^2} + H'''(\xi) \cdot \frac{(1-x)^3}{(N+3)^3}$$

$$H\left(\frac{(N+2)x}{N+3}\right) - H(x) = H'(x) \cdot \frac{-x}{N+3} + \frac{1}{2}H''(x) \cdot \frac{x^2}{(N+3)^2} + H'''(\tilde{\xi}) \cdot \frac{-x^3}{(N+3)^3}.$$

Thus

$$\mathfrak{J}_H(X) = \frac{x(1-x)}{2(N+3)^2}H''(x) + \left(\frac{x(1-x)^3}{(N+3)^3}H'''(\xi) - \frac{x^3(1-x)}{(N+3)^3}H'''(\tilde{\xi})\right). \tag{36}$$

Because H is smooth and $\left[\frac{(N+2)x}{N+3}, \frac{(N+2)x+1}{N+3}\right]$ is compact, the third derivative of H is bounded. Therefore, as N increases,

$$\mathfrak{J}_H(X) = \mathfrak{D}_{\text{beta}} H(x) + O\left(\frac{1}{(N+3)^3}\right).$$

Moreover, for log scoring rule H_{ln} , its second derivative is $\frac{1}{\ln(2)x(1-x)}$ and the third derivative is $\frac{-1}{x^2 \ln 2} + \frac{1}{(1-x)^2 \ln 2}$, by Eq. (36), the difference $2(\ln 2)(N+3)^2 \mathfrak{J}_{H_{\text{ln}}}(X) - 1$ is o(1) for all $\delta_N \leq x \leq 1 - \delta_N$. Formally

$$\begin{split} &|2(\ln 2)(N+3)^2\,\mathfrak{J}_{H_{\ln}}(X)-1|\\ &=\frac{2}{N+3}\left|x(1-x)^3\left(\frac{-1}{\xi^2}+\frac{1}{(1-\xi)^2}\right)-x^3(1-x)\left(\frac{-1}{\tilde{\xi}^2}+\frac{1}{(1-\tilde{\xi})^2}\right)\right|\\ &\leq \frac{2}{N+3}\left(\frac{x}{\xi^2}+\frac{(1-x)^3}{(1-\xi)^2}+\frac{x^3}{\tilde{\xi}^2}+\frac{1-x}{(1-\tilde{\xi})^2}\right)\\ &\leq \frac{2}{N+3}\left(\frac{1}{x}+\frac{(1-x)^3}{(1-\xi)^2}+\frac{x^3}{\tilde{\xi}^2}+\frac{1}{1-x}\right)\\ &\leq \frac{2}{N+3}\left(\frac{1}{x}+\frac{N+3}{N+2}(1-x)+\frac{N+3}{N+2}x+\frac{1}{1-x}\right)\\ &\leq \frac{2}{N+3}\left(\frac{1}{x}+\frac{N+3}{N+2}(1-x)+\frac{N+3}{N+2}x+\frac{1}{1-x}\right)\\ &\leq \frac{4}{(N+3)\delta_N}+\frac{2}{N+2} \end{split} \qquad (\xi\leq \frac{(N+2)x+1}{N+3} \text{ and } \frac{(N+2)x}{N+3}\leq \tilde{\xi})\\ &\leq \frac{4}{(N+3)\delta_N}+\frac{2}{N+2} \end{split}$$

Therefore, $\text{Obj}_{\bar{\mathcal{X}}_{Beta}^{N}}(H_{\ln}) = \frac{1}{2(\ln 2)(N+3)^{2}} + o(\frac{1}{(N+3)^{2}}).$

Proof of Lemma 5.16. The proof is very similar to Theorem 3.5 in Evans [11]. Informally although the matrix A is only positive definition in a linear subspace, we also only need to ensure the maximum position in the same space. Because v_1, \ldots, v_k is an orthonormal basis of \mathcal{V} , we can add d-k vectors v_{k+1}, \ldots, v_d such that v_1, \ldots, v_d is an orthonormal basis for \mathbb{R}^d .

Suppose x^* is an interior maximum point of h in W. Since h is smooth, $\nabla^2 h(x^*)$ is symmetric matrix, and $v_\ell^\top \nabla^2 h(x^*) v_\ell \leq 0$ for all $\ell \in [k]$ by the second derivative test. Therefore, there exists $\lambda_1, \ldots, \lambda_d$ such that $\lambda_1, \ldots, \lambda_k \leq 0$, and $\nabla^2 h(x^*) = \sum_\ell \lambda_\ell v_\ell v_\ell^\top$.

First we assume u is a strict supersolution such that $0 < \mathcal{L}h(x)$ for all $x \in W$. For $x = x^*$, we have

$$0 < \mathcal{L}h(x^*)$$

$$= \sum_{i,j} A_{ij}(x^*) \partial_i \partial_j h(x^*)$$

$$= \sum_{i,j} A_{ij}(x^*) (\nabla^2 h(x^*))_{ij}$$

$$= \operatorname{Tr} \left(A(x^*) (\nabla^2 h(x^*))^\top \right) \qquad (\sum B_{ij} C_{ij} = \operatorname{Tr} (BC^\top))$$

$$= \operatorname{Tr} \left(A(x^*) (\nabla^2 h(x^*)) \right) \qquad (\nabla^2 u \text{ is symmetric})$$

$$= \operatorname{Tr} \left(A(x^*) (\sum_{\ell} \lambda_{\ell} v_{\ell} v_{\ell}^\top) \right) \qquad (\nabla^2 h(x^*) = \sum_{\ell} \lambda_{\ell} v_{\ell} v_{\ell}^\top)$$

$$= \sum_{\ell=1}^d \lambda_{\ell} \operatorname{Tr} \left(v_{\ell}^\top A(x^*) v_{\ell} \right).$$

Because $null(A(x^*)) = \mathcal{V}^{\perp}$, $A(x^*)v_{\ell} = 0$ for all $k < \ell \le d$. So the $\sum_{\ell=1}^{d} \lambda_{\ell} \operatorname{Tr} \left(v_{\ell}^{\top} A(x^*) v_{\ell} \right) = \sum_{\ell=1}^{k} \lambda_{\ell} \operatorname{Tr} \left(v_{\ell}^{\top} A(x^*) v_{\ell} \right)$. However, since for all $1 \le \ell \le k$ $\lambda_{\ell} < 0$ and $v_{\ell}^{\top} A(x^*) v_{\ell} > 0$, we have

$$\mathcal{L}h(x^*) = \sum_{\ell=1}^d \lambda_\ell \operatorname{Tr} \left(v_\ell^\top A(x^*) v_\ell \right) \le 0.$$

Thus, we reach a contradiction, so x^* cannot be an interior maximum point. Furthermore $\max_{cl(W)} h = \max_{bd(W)} h$ because h is continuous on cl(W).

However, if h only satisfies $\mathcal{L}h \geq 0$, we define $h_{\epsilon}(x) := h(x) + \epsilon \exp(\langle v_1, x \rangle)$ for $\epsilon > 0$. Thus, $\mathcal{L}h_{\epsilon}(x) = \mathcal{L}h(x) + \epsilon v_1^{\top} A(x) v_1 \exp(\langle v_1, x \rangle) > 0$. Then apply the above argument we have $\max_{cl(W)} h_{\epsilon} = \max_{bd(W)} h_{\epsilon}$. By taking ϵ to zero, we have $\max_{cl(W)} h = \max_{bd(W)} h$.

Finally, suppose h has an interior maximum and not constant. Set $Y:=\{y:h(y)=\max h\}$ and $Z=W\setminus Y$. We choose a point $z\in Z$ such that dist(z,Y)< dist(z,bd(W)), and B denote the largest ball that center z whose interior lies in Z. Then there exists some point $y^*\in Y$ with $y^*\in bd(B)$. Thus, Z satisfies the interior ball condition at y^* . By Hopf's Lemma [11], we have $\frac{\partial h}{\partial \nu}(y^*)>0$ where ν is the vector normal to B at y^* . But this is a contradiction, since h attains its maximum at y^* , we have $\nabla h(y^*)=0$. \square