PIER Working Paper

18-008

# Overabundant Information and Learning Traps

ANNIE LIANG
University of Pennsylvania
Department of Economics

XIAOSHENG MU
Harvard University

March 27, 2018

# Overabundant Information and Learning Traps*

Annie Liang[†]     Xiaosheng Mu[‡]

March 27, 2018

**Abstract**

We develop a model of learning from overabundant information: Agents have access to many sources of information, where observation of all sources is not necessary in order to learn the payoff-relevant unknown. Short-lived agents sequentially choose to acquire a signal realization from the best source for them. All signal realizations are public. Our main results characterize two starkly different possible long-run outcomes, and the conditions under which each obtains: (1) efficient information aggregation, where signal acquisitions eventually achieve the highest possible speed of learning; (2) "learning traps," where the community gets stuck using an suboptimal set of sources and learns inefficiently slowly. A simple property of the correlation structure separates these two possibilities. In both regimes, we characterize which sources are observed in the long run and how often.

## 1 Introduction

Most informational environments are environments of *informational overabundance*: there are more sources of information than any individual can attend to. A key decision that individuals must make in such settings is which sources to observe. These informational choices can have significant consequences: when agents choose information at different times, the informational choices of earlier agents impose externalities

on those of later agents. We demonstrate the possibility for a small number of early inefficient choices to propagate across time, and characterize the exact conditions on the available informational sources under which this can happen.

To fix ideas, consider a development economist working on the question of how access to microfinance affects poverty reduction. To contribute towards our understanding of these questions, the development economist may run a randomized control trial in which microcredit is provided to selected individuals. The key constraint is that the development economist cannot choose to run this experiment in an ideal, "representative," sample of the global population, but must choose a single setting: a village in the suburbs of Bangalore, for example. Experiments conducted in any individual location provide an inevitably *biased* estimate of the effect of microcredit on poverty, although the exact extent or direction of this bias may be unknown. Moreover, the outcomes of experiments conducted in different locations are related by a potentially complex pattern of correlations: two different villages in the suburbs of Bangalore may yield close estimates, whereas the same experiment in Morocco may yield a rather different outcome.

From the perspective of society, there is an optimal allocation of experiments across locations, but the incentives provided for any single development economist can be in conflict with this goal. In particular, each development economist may care only about the *marginal* contribution of her experiment towards understanding (for example, if this is the relevant criterion for publication). The question of interest is how the choices of these short-lived decision-makers influence the choices of subsequent decision-makers, and the speed at which the community learns.

We study these questions within a sequential learning model. There is an unknown payoff-relevant state $\omega$ (the true impact of microcredit on poverty) and additionally $K-1$ possible *biases* or *confounding terms*, labeled $b_1, \ldots, b_{K-1}$. Information sources (locations where RCTs can be implemented) are modeled as different linear combinations of these $K$ unknowns, plus an independent Gaussian error. For example, it may be that microcredit is a relatively more effective policy tool in India (relative to other countries), and there may additionally be regional variation across India regarding suitability of this intervention. Thus, running a RCT in a particular Indian village results in an estimate that is affected by a country-specific and also a region-specific bias. We will say that there are "overabundant" sources of information if experiments need not be conducted at every possible location in order to recover $\omega$.

Agents (development economists), indexed by $t \in \mathbb{N}$, sequentially choose locations

at which to conduct their experiment. Each RCT generates a *public* independent realization of the corresponding (biased) signal of $\omega$. The feature that information is public contrasts with the classic sequential learning model (Banerjee, 1992; Bikhchandani, Hirshleifer and Welch, 1992; Smith and Sorenson, 2000). Public information permits us to focus on the externalities created by choices of kind of information, as opposed to the more frequently studied frictions that emerge from problems of inference. We assume that each agent takes an action based on all information acquired up to (and including) himself in order to maximize a private objective that depends only on $\omega$. As we will discuss, the specific choice of payoff function is not important: the information choice that maximizes payoffs for each agent will be the one that maximizes that period's reduction of (society's) uncertainty about $\omega$.

We show that there are two (exhaustive) possible long-run outcomes:

- *Efficient information aggregation:* when sources are related in a particular way (that we characterize), the pattern of experimentation eventually mimics the "best" sampling pattern across locations. In this case, development economists eventually conduct RCTs only at the set of locations which jointly maximize the speed of society's learning about $\omega$.

- *Learning traps:* otherwise, depending on the common prior, there is a set of possible long-run observation sets (that we characterize), including ones in which agents become "trapped" conducting RCTs at a set of locations which do reveal $\omega$, but do so inefficiently slowly.

Our main contribution is to demonstrate that which of these outcomes obtains depends on a simple property of the correlation structure across sources.

Formally, we refer to any set of informational sources that reveal $\omega$ as a *spanning set*, and say that a spanning set is *minimal* if no proper subset of its sources reveals $\omega$. Whether the long-run outcome is efficient aggregation or a learning trap depends critically on whether there exists a minimal spanning set that is of lower-dimension than the state space. The key intuition refers back to an observation used in Sethi and Yildiz (2016); recall that an agent who observes a biased source learns *both* about the payoff-relevant state and also about the source's own bias. In our setting, where biases are correlated across sources, there is a further spillover effect: learning from a biased source helps agents to understand the biases of all sources that are correlated with it. Suppose now that agents repeatedly observe RCT outcomes at a set of $K$ locations that collectively reveal all $K-1$ unknown biases in addition to the state $\omega$ (say that

these sources have "full rank"). Then, every time an agent conducts a new RCT at one of these locations, he improves society's understanding about how to interpret RCT estimates not only at these locations, but also at all other locations. It can be shown that eventually agents come to evaluate locations by an objective asymptotic criterion that is prior-independent. We thus present the following positive result: if every minimal spanning set has full rank, then long-run acquisitions are optimal, independently of the prior belief.

In contrast, if it is possible to learn $\omega$ by conducting RCTs at fewer than $K$ locations, then inefficient long-run learning may obtain. Intuitively, RCTs at $k < K$ locations provide bounded positive spillovers for sources outside of the set. This is because agents can at most learn $k - 1$ unknown biases from these sources, while the other sources may depend on the remaining $K - k$ biases. Thus, the community's understanding of locations where experiments have not yet been conducted need not improve. Formally, we show that for every minimal spanning set that is "best" in its subspace, there is an open set of priors given which this set of locations is exclusively observed in the long run. The implied inefficiency—measured as the ratio of the optimal speed of learning and the achieved speed of learning—can be an arbitrarily large constant.

Our work combines ideas from two literatures. First, recent work (Sethi and Yildiz, 2016; Che and Mierendorff, 2017; Fudenberg, Strack and Strzalecki, 2017; Liang, Mu and Syrgkanis, 2017; Mayskaya, 2017; Sethi and Yildiz, 2017) studies choice of information from a finite set of information sources. We build specifically upon Liang, Mu and Syrgkanis (2017), which introduced the framework we describe in Section 3 under a restriction that the number of sources and states are the same (thus ruling out the possibility of informational overabundance, which is the focus here). Our work also builds on Sethi and Yildiz (2016, 2017), which study long-run acquisitions from a large number of Gaussian sources. Our model differs from this related work in a few key ways: First, Sethi and Yildiz (2016, 2017) consider stochastic error variances, so that the "best" sources vary from period to period, while we fix noise variances, so that there is (generically) a unique "best" asymptotic set. Second, Sethi and Yildiz (2016, 2017) focus on correlation structures that fall under our Theorem 2, for which long-run acquisitions do not necessarily achieve efficient learning, while we explore also those correlation structures that lead to optimal learning. Thus, the welfare comparisons that we make here are particular to our framework.

Finally, our model relates to the social learning and herding literatures (Banerjee,

1992; Bikhchandani, Hirshleifer and Welch, 1992), which consider information aggregation by short-lived agents who sequentially acquire information. At a high level, the externality identified in our paper is related to the classic externality from this literature: in both settings, the precision of public information can grow inefficiently slowly because of endogenous information acquisitions driven by past choices. But in the present paper, all signal realizations are publicly and perfectly observed, which turns off the inference problem essential to the existence of cascades in standard herding models. Our focus is on a new mechanism, in which externalities arise through choice of *kind* of information; as we will discuss, this externality has a rather different nature structure from those studied previously.

Finally, recent papers introducing costly information acquisition to the sequential learning model include Burguet and Vives (2000), Ali (2017), and Mueller-Frank and Pai (2016). Relative to this work, our paper considers choice from a fixed set of information sources (with a capacity constraint), in contrast to choice from a flexible set of information sources (with a cost on precision).

# 2   Example

Development economists sequentially run RCTs to uncover an unknown parameter $\omega \in \mathbb{R}$; for example, the impact of microcredit on poverty. Each RCT yields a noisy estimate of $\omega$ that is biased by the specific social and market environment in which it are conducted. The nature of these biases is correlated across different environments, and in particular, studies that are conducted in similar environments are biased in similar ways.

At each period $t = 1, 2, \ldots$, an economist chooses to run a RCT at the location that will reduce (societal) uncertainty about $\omega$ as much as possible *in that period*. The results of all experiments are public. The question of interest is whether long-run acquisitions will efficiently aggregate information about $\omega$, and whether RCTs will eventually be conducted at the "best" set of sites. Below we contrast two patterns of correlations across potential RCT sites.

**Setting 1:** RCTs can be conducted in Morocco or India, and there is additionally heterogeneity over how urban the location is. Write $b_{\text{Morocco}}$ for a Morocco-specific bias, $b_{\text{India}}$ for an India-specific bias, and $b_{\text{Urban}}$ for an urban bias. Additionally, suppose that there is heterogeneity in how precise the estimates are across sites. The

following is an example set of sources with these biases:

$$X_1 = \omega + b_{\text{India}} + b_{\text{Urban}} + \epsilon_1$$

$$X_2 = 2\omega + b_{\text{India}} + 2b_{\text{Urban}} + \epsilon_2$$

$$X_3 = \omega + b_{\text{India}} + \epsilon_3$$

$$X_4 = 2\omega + b_{\text{Morocco}} + b_{\text{Urban}} + \epsilon_4$$

$$X_5 = 5\omega + b_{\text{Morocco}} + 2b_{\text{Urban}} + \epsilon_5$$

$$X_6 = 3\omega + b_{\text{Morocco}} + \epsilon_6$$

where error terms are i.i.d standard normals.[1]

Given this correlation structure, application of our subsequent Theorem 1 yields that the community's speed of learning is maximized if agents eventually sample only from the sites in Morocco (corresponding to observation of $\{X_4, X_5, X_6\}$). But there is a set of priors given which agents exclusively run RCTs in India in the long run (this follows from our subsequent Theorem 2). Intuitively, observations of $X_1$, $X_2$, and $X_3$ provide information not only about the parameter of interest $\omega$, but also about the India-specific bias $b_{\text{India}}$. Thus, each RCT that is conducted in India helps future economists to de-bias subsequent RCTs in India. This can create a self-reinforcing sequence of choices, where observations of $X_1$, $X_2$, and $X_3$ increase the value of future observations of these sources relative to the sources $X_4$, $X_5$, and $X_6$.

In contrast, consider the following setting:

**Setting 2:** There is heterogeneity across sites along the following dimensions: whether lending is restricted to women, whether the location is urban, and the intensity to which lending is targeting towards entrepreneurs. Write $b_{\text{WomenOnly}}$, $b_{\text{Entrepeneurs}}$, and $b_{\text{Urban}}$ for these biases. As before, there is additionally heterogeneity in the precision of estimates across sites. The following is an example set of sources with these

---

[1]Note that larger coefficients on $\omega$ correspond to more precise estimates.

biases:

$$X_1 = \omega + b_{\text{WomenOnly}} + 2b_{\text{Entrepeneurs}} + b_{\text{Urban}} + \epsilon_1$$

$$X_2 = 2\omega + b_{\text{WomenOnly}} + b_{\text{Urban}} + \epsilon_2$$

$$X_3 = \omega + b_{\text{Entrepeneurs}} + \epsilon_3$$

$$X_4 = 2\omega + b_{\text{WomenOnly}} + b_{\text{Entrepeneurs}} + b_{\text{Urban}} + \epsilon_4$$

$$X_5 = 5\omega + b_{\text{Entrepeneurs}} + b_{\text{Urban}} + \epsilon_5$$

$$X_6 = 3\omega + b_{\text{WomenOnly}} + 3b_{\text{Entrepeneurs}} + \epsilon_6$$

where error terms are i.i.d standard normals.

Again applying our subsequent Theorem 1, the community's speed of learning is maximized if agents eventually sample only from sites $X_1$, $X_2$, $X_5$, and $X_6$. In contrast to the above setting, however, economists will necessarily eventually conduct studies only in the best subset of sites, and their information acquisitions will approximate the optimal sampling over these sites. This conclusion holds irrespective of the common prior.

What differentiates these two informational environments? The framework that we present in the next section includes these two correlation structures above as special cases, and our subsequent analysis demonstrates that a simple property of the correlations (across sources) separates the two settings.

# 3 Framework

## 3.1 Model

There is a persistent unknown payoff-relevant state $\omega$, and additionally there are $K-1$ persistent unknown biases $b_1, \ldots, b_{K-1}$.[2] Throughout, it will be convenient to write $\theta = (\omega, b_1, \ldots, b_{K-1})'$ for the $K$-dimensional vector of unknowns.[3] Assume that $\theta$ follows a multivariate normal distribution $\mathcal{N}(\mu^0, V^0)$, where the prior covariance matrix $V^0$ has full rank.[4] Agents have access to $N$ different sources of information,

---

[2] We refer to the biases also as "states", but we will say that $\omega$ is the only "payoff-relevant state" among the $K$ states.

[3] All vectors in this paper are column vectors.

[4] This assumption rules out linear dependence across the state and biases. If there is linear dependence, we may work with a smaller set of state and biases without changing the model.

and observation of source $i$ corresponds to an independent realization of the signal

$$X_i^t = \langle c_i, \theta \rangle + \epsilon_i^t, \quad \epsilon_i^t \sim \mathcal{N}(0, 1).$$

where each $c_i = (c_{i1}, \ldots, c_{iK})'$ is a constant vector, and the noise terms $\epsilon_i^t$ are independent from each other and over time. Our assumption that noise terms have unit variance is without loss since the coefficients $c_i$ are unrestricted. We let $C$ denote the $N \times K$ matrix whose $i$-th row is $c_i'$.

A countably infinite number of agents, indexed by $t \in \mathbb{N}$, move sequentially. Each agent $t$ acquires an independent realization of one of the $N$ signals, and then chooses an action $a \in A$ to maximize an individual objective $u_t(a, \omega)$. He bases his action on the realization of his own signal acquisition, as well as the history of signal acquisitions and realizations thus far. (Thus, all signal realizations are public.)

Payoff functions may differ across agents, but we assume that the decision problems are non-trivial in the following way.

**Assumption 1** (Payoff Sensitivity to Mean)**.** *For every $t$, any variance $\sigma^2 > 0$ and any action $a^* \in A$, there exists a positive measure of $\mu_1$ for which $a^*$ does not maximize $\mathbb{E}[u_t(a, \omega) \mid \omega \sim \mathcal{N}(\mu_1, \sigma^2)]$.*

In words, holding the belief variance fixed, the expected value of $\omega$ affects the optimal action to take. A sufficient condition for Assumption 1 is that for every agent $t$ and every action $a^*$, there exists some other action $\hat{a}$ such that $u_t(\hat{a}, \omega) > u_t(a^*, \omega)$ as $\omega \to +\infty$ or as $\omega \to -\infty$. That is, we require that the two limiting states $\omega \to +\infty$ and $\omega \to -\infty$ yield different optimal actions. This is true for all natural applications.

We now introduce some terminology to describe possible information environments. Let $[N] = \{1, \ldots, N\}$ index the set of signals. We will call a set of signals $\mathcal{S} \subset [N]$ a *spanning set* if the vectors $\{c_i : i \in \mathcal{S}\}$ span the coordinate vector $e_1 = (1, 0, \ldots, 0)'$, so that it is possible to learn $\omega$ by exclusively observing signals from $\mathcal{S}$. If $\mathcal{S}$ is spanning, and no proper subset of $\mathcal{S}$ is spanning, then we will refer to $\mathcal{S}$ as a *minimal spanning set*.

Throughout this paper, we assume that the complete set of signals $[N]$ is spanning, so that $\omega$ can be recovered by observing all signals infinitely often.[5] This assumption

---

[5]Our results do extend to situations where $\omega$ is *not* identified from the $N$ available signals. To see this, first note that we can always take a linear transformation and work with the following equivalent model: the state vector $\tilde{\theta}$ is $K$-dimensional *standard Gaussian*, each signal $X_i$ is $\tilde{c}_i'\tilde{\theta} + \epsilon_i$ and the payoff-relevant parameter is $u'\tilde{\theta}$ for some fixed vector $u$. Let $V$ be the subspace of $\mathbb{R}^K$

nests two interesting cases. Say that the informational environment has *exactly sufficient information* if $[N]$ is minimally spanning. Then, it is possible to recover $\omega$ by observing each informational source infinitely often, but not by observing any proper subset of sources.

Our main setting of interest, which we refer to as *informational overabundance*, obtains when $[N]$ is spanning but not minimally spanning. Then, there are multiple different sets of signals which allow for recovery of $\omega$, and a key point of our analysis is to compare the set of sources that "should" be observed in the long run and the set of sources that is in fact observed in the long run. Except for trivial cases, informational overabundance corresponds to $N > K$ (more signals than states).[6]

## 3.2 Interpretation

As mentioned in the introduction, the framework can be used to describe the following setting: development economists sequentially choose settings under which to run RCTs, each of which provides a biased estimate of a parameter of interest (e.g. the effect of microcredit on poverty). Crucially, biases may be flexibly correlated across these choices, and may reflect a composition of many different kinds of biases. For example, the RCT estimates may be biased by the country in which the RCT is conducted, the size of the lending group, as well as various restrictions on lending. Economists choose settings to maximize a personal and myopic objective: for example, they may desire for their experiment to reduce (societal) uncertainty about the unknown parameter as much as possible, and thus choose the location that allows for the greatest immediate reduction in posterior variance about the unknown parameter. We will subsequently describe our main results relative to this primary interpretation.

---

spanned by $\tilde{c}_1, \ldots, \tilde{c}_N$. Then consider the projection of $u$ onto $V$: $u = v + w$ with $v \in V$ and $w$ orthogonal to $V$. This enables us to write $u'\tilde{\theta} = v'\tilde{\theta} + w'\tilde{\theta}$. By assumption, the random variable $w'\tilde{\theta}$ is independent from any random variable $c'\tilde{\theta}$ with $c \in V$ (because they have zero covariance). Thus the uncertainty about $w'\tilde{\theta}$ is not reduced upon any signal observation. Consequently the agents only seek to learn about $v'\tilde{\theta}$, returning to the case where the payoff-relevant parameter *is* identified.

[6]It is possible for $\omega$ to be "overidentified" from a set of $N \leq K$ signals, e.g. $X_1 = \omega + \epsilon_1$, $X_2 = \omega + b_1 + b_2 + \epsilon_2$, and $X_3 = b_1 + b_2 + \epsilon_3$. In this case, the set $\{X_1, X_2, X_3\}$ is spanning, but not minimally spanning since both of its subsets $\{X_1\}$ and $\{X_2, X_3\}$ are also spanning. Although $N = K = 3$ in this example, it is equivalent to a model in which there is a single bias $\tilde{b}_1 = b_1 + b_2$, and the three signals are rewritten $X_1 = \omega + \epsilon_1$, $X_2 = \omega + \tilde{b}_1 + \epsilon_2$ and $X_3 = \tilde{b}_1 + \epsilon_3$. Then, we do have $N > K$ in this alternative model.

The basic framework applies also, however, to a variety of other settings. We briefly describe a few below.

*Biased media sources and experts.* A related interpretation takes each information source to be a media outlet or an expert, where biases are correlated across outlets, and potentially multi-dimensional. For example, a newspaper may simultaneously want to bias its reporting towards its liberal readership, its conservative donors, the values of its editorial board, and also in a way that increases sales. Different weights on these objectives across media sources will result in different (but correlated) slants on the same issue.

*Culture:* Interpret each source as feedback from an individual or a social group, where accurate information processing requires an understanding of the individual's (or group's) cultural context. For example, a computer scientist presenting for the first time at an economics seminar may not know how to interpret aggressive questions from the audience, which could indicate skepticism, or simply a high level of interest and a more aggressive cultural norm. The first interaction with a given social group thus yields noisy feedback, but repeated interactions allow one to develop a better understanding for that source. The biases in our framework can be interpreted as a reduced form model for "communication noise" that reduces with repeated observation.

*Attribute sampling:* Suppose that a consumer good is described by $K$ attributes $\tilde{\theta}_1, \ldots, \tilde{\theta}_K$, and its unknown quality $\omega = \sum_{k=1}^{K} \alpha_k \tilde{\theta}_k$ is a linear combination of its qualities along each of these dimensions. For example, consumers want to learn about the quality of a new laptop, where quality is a linear combination of a large number of attributes, including ease of use, battery life, and screen resolution. Agents can learn about each of these attributes through different kinds of inspections, which provide noisy information about different linear combination of attributes. This model can be rewritten in our framework above, where the state vector is $(\omega, b_1, \ldots, b_{K-1})$ and each "bias" term $b_i = \tilde{\theta}_i$ for $1 \leq i \leq K - 1$.

# 4 Preliminaries

Each agent $t$ faces a *history* $h^{t-1} \in ([N] \times \mathbb{R})^{t-1} = H^{t-1}$ consisting of all past signal choices and realized signal values. A *strategy* for agent $t$ is a measurable map from all $(t-1)$-length histories to signals—that is, $S : H^{t-1} \to [N]$, where $S(h^{t-1})$ represents

the signal choice in period $t$ following history $h^{t-1}$.[7]

Each agent $t$'s marginal belief about $\omega$ (updated to his own signal acquisition and all information revealed by past agents) is Gaussian, so we can write $\omega \sim \mathcal{N}(\mu_1^t, V_{11}^t)$ for this belief. His maximum expected payoff is

$$\max_{a \in A} \mathbb{E}[u_t(a, \omega) \mid \omega \sim \mathcal{N}(\mu_1^t, V_{11}^t)] \tag{1}$$

Each agent $t$ chooses the signal that maximizes (1). Observe further that since beliefs are Gaussian, the agent's posterior variance $V_{11}$ about $\omega$ following $q_i$ observations of each signal $i$ can be written as a deterministic function $V_{11} = f(q_1, \ldots, q_N)$. In particular, the posterior variance does not depend on signal realizations; see Appendix B for the complete (closed-form) expression.

We use the following lemma, which says that the signal choice that achieves the greatest reduction in posterior variance also maximizes expected payoffs:

**Lemma 1** (Liang, Mu and Syrgkanis (2017)). *The optimal signal acquisition for every agent, at every history, is the signal that minimizes current posterior variance about $\omega$.*

Using this lemma, we can track society's acquisitions in the following way. Write $m(t) = (m_1(t), \ldots, m_N(t))$ for the *division* over signals at time $t$, where $m_i(t)$ is the number of times that signal $i$ has been observed up to and including time $t$. Then, $m(t)$ evolves deterministically according to the following rule: $m(0)$ is the zero vector, and for each $t \geq 0$,

$$m_i(t+1) = \begin{cases} m_i(t) + 1 & \text{if } f(m_i(t) + 1, m_{-i}(t)) \leq f(m_j(t) + 1, m_{-j}(t)) \; \forall j. \\ m_i(t) & \text{otherwise.} \end{cases}$$

That is, in each period $t$ the division vector increases by 1 in exactly one coordinate, corresponding to the signal that allows for the greatest immediate reduction in posterior variance.[8]

We are primarily interested in the *long-run* signal acquisitions: which sources are observed in the long run, and how often? Below, we characterize the *asymptotic frequency* $\lim_{t \to \infty} m_i(t)/t$ with which source $i$ is observed, and discuss the *asymptotic observation set*, meaning the signals that are observed with positive frequency in the long-run. Our subsequent results in Section 6 show that these limits are well-defined.

---

[7]Since information is public, agents do not need to additionally condition on past actions.

[8]We allow ties to be broken arbitrarily, so there may be multiple paths $m(t)$.

We will also be interested in how these long-run signal acquisitions compare to "optimal" acquisitions that a social planner might impose. Below, we begin by characterizing an optimal benchmark, corresponding to the choices that permit the greatest amount of information revelation (Section 5) and then turn to the community's long run acquisitions.

# 5  Benchmark: Social Planner

Throughout we evaluate society's acquisitions relative to an "optimal" benchmark, which we define as the limit of a sequence of solutions to finite horizon problems. Consider a social planner who takes an action $a \in A_{SP}$ on behalf of the society at some large period $t$, with payoff function $u_t^{SP}(a, \omega)$.

The social planner's payoffs are maximized if the history of $t$ signal acquisitions are allocated across signals in the following way (see Lemma 3 in Liang, Mu and Syrgkanis (2017)):

$$n(t) \in \underset{(q_1,\ldots,q_K):q_i\in\mathbb{Z}^+,\sum_i q_i=t}{\operatorname{argmin}} f(q_1,\ldots,q_K).$$

That is, any allocation of the $t$ observations that minimizes posterior variance about $\omega$ will maximize the social planner's payoffs.[9] Generically, there is a unique optimal division vector $n(t)$ for every $t$.[10]

We can interpret each $n(t)$ as the optimal social benchmark for the finite horizon problem with final period $t$. The limiting frequencies $\lim_{t\to\infty} n(t)/t$ are well-defined under a subsequent condition, and we refer to these as the *optimal long-run frequencies*. If agents repeatedly choose signals according to these limiting frequencies, then the empirical distribution over signals at any large $t$ will be arbitrarily close to $n(t)/t$. Since payoffs are continuous in signal frequencies, this stationary rule also approximates aggregate payoffs under a $\delta$-discounting criterion when $t$ is sufficiently large.[11]

---

[9]In more detail, suppose the planner can dictate signal choices to maximize the expected payoff after $t$ periods. Then she should use any strategy that observes each signal $i$ exactly $n_i(t)$ times. In particular, there exists an optimal strategy that does not condition on signal realizations.

[10]Throughout the paper, "generic" means with probability 1 for signal coefficients $c_{ij}$ randomly drawn from a full support distribution on $\mathbb{R}^{NK}$.

[11]We conjecture that the limiting frequencies $n(t)/t$ also "maximize" the $\delta$-discounted objective when agents share the same utility function. Formally, for any fixed $\delta$, let $\mathcal{S}$ denote a strategy that

Below, we break up the characterization of $\lim_{t \to \infty} n(t)/t$ into two cases: in Section 5.1 we discuss settings with *exactly sufficient information*, where all sources must be observed infinitely often to recover $\omega$. In Section 5.2 we consider *informational overabundance*, where asymptotic learning can occur from a strict subset of sources.

## 5.1 Exactly Sufficient Information

Settings of exactly sufficient information, in which all signals must be observed in order to recover $\omega$, have fewer signals than states ($N \leq K$).[12] Moreover, in such settings, it is possible to decompose the first coordinate vector as a (unique) linear combination of all of the available signals: $e_1 = \sum_{i=1}^{N} \beta_i \cdot c_i$, where the coefficients $\beta_i$ are nonzero real numbers.

Assume for now that $N = K$. Then each $\beta_i$ admits the simple form $\beta_i = |[C^{-1}]_{1i}|$.[13] In our prior work, we showed that the optimal frequency with which each signal $i$ is observed asymptotically is proportional to its coefficient $\beta_i$:

**Proposition 1** (Liang, Mu and Syrgkanis (2017)). *Suppose $N = K$ and there is exactly sufficient information. Then for every signal $i$, the optimal count satisfies*

$$n_i(t) = \frac{|\beta_i|}{\sum_{j=1}^{N} |\beta_j|} \cdot t + O(1). \tag{2}$$

Throughout, $O(1)$ represents an error term that remains bounded as $t \to \infty$.

The above result extends to environments with strictly fewer signals than states ($N < K$), under the assumption of exactly sufficient information. This follows from the observation that any environment with $N < K$ can be transformed into an equivalent environment with $N = K$. For example, suppose the available signals are $X_1 = \omega + b_1 + b_2$ and $X_2 = \omega - b_1 - b_2$, so that the number of states ($K = 3$) exceeds

---

maximizes $\sum_{t=1}^{\infty} \delta^t \cdot u(a_t, \omega)$. Further let $d^\delta(t)$ be the vector of signal counts at time $t$, under strategy $\mathcal{S}$. Then we conjecture that for all $\delta$ close to 1, $\lim_{t \to \infty} d^\delta(t)/t = \lim_{t \to \infty} n(t)/t$.

[12]It is a simple fact in linear algebra that whenever the $K$-dimensional vector $e_1$ is spanned by vectors $c_1, \ldots, c_N$, it must be spanned by some $K$ of these vectors.

[13]These coefficients can be interpreted in the following way: Suppose that a single realization of each signal $X_i$ is observed. The random vector describing these realizations can be written $Y = C\theta + \varepsilon$, where $\varepsilon$ is the $K \times 1$ vector of error terms. Given a realization $Y$, the best linear unbiased estimate for the payoff-relevant state $\omega$ is then $\hat{\omega} = [C^{-1}Y]_{11}$. If we perturb the vector $Y$ by $\eta$ in coordinate $i$, the estimate for $\omega$ changes (in magnitude) by $\beta_i \cdot \eta = |[C^{-1}]_{1i}| \cdot \eta$. Thus, the larger $\beta_i$ is, the more the estimate $\hat{\omega}$ responds to changes in the realization of signal $X_i$. Informally, the larger $\beta_i$ is, the "more important" signal $X_i$ is in determining the estimate of $\omega$.

the number of signals ($N = 2$). We can define a new state $\tilde{b}_1 = b_1 + b_2$ and rewrite $X_1 = \omega + \tilde{b}_1$ and $X_2 = \omega - \tilde{b}_1$. In this equivalent model, the number of states and signals are the same ($N = K = 2$).

This transformation applies in general: we can always choose $N$ states (including $\omega$), each a linear combination of the original $K$ states, and re-define the original signals to be linear combinations of the new states. The transformed model is equivalent to the original problem, but it satisfies the conditions of Proposition 1. Thus, dropping the requirement that $N = K$, we obtain the following corollary:

**Corollary 1.** *Suppose there is exactly sufficient information. Then each $n_i(t) = \frac{|\beta_i|}{\sum_{j=1}^{N} |\beta_j|} \cdot t + O(1)$, for $1 \leq i \leq N$.*

Thus, a sampling procedure in which each signal $i$ is chosen with probability $|\beta_i| / \sum_{j=1}^{N} |\beta_j|$ each period will eventually approximate $n(t)$. The subsequent corollary describes the speed of learning along the sequence $n(t)$.

**Corollary 2.** *The minimum posterior variance after $t$ observations satisfies the following approximation:*

$$f(n(t)) = \min_{\sum_{i=1}^{N} q_i = t} f(q_1, \ldots, q_N) \sim \left( \sum_{i=1}^{N} |\beta_i| \right)^2 / t.$$

*where the notation "$F(t) \sim G(t)$" means $\lim_{t \to \infty} \frac{F(t)}{G(t)} = 1$.*

Thus, eventually agents approximate a posterior variance of $\left( \sum_{i=1}^{N} |\beta_i| \right)^2 / t$.

A key property of the result above is that every signal is viewed with positive frequency in the long run. This is natural when all signals must be observed in order to recover $\omega$. But when there is informational overabundance, a question emerges regarding *how many* and *which* signals will be observed asymptotically.

## 5.2 Informational Overabundance

Suppose now that the number of signals exceeds the number of states ($N > K$). In *informationally overabundant* environments, $\omega$ can be learned by exclusively observing signals from any of several distinct spanning sets. In principle, speed of learning can be further improved by combining observations from multiple spanning sets. We put this aside for the moment, and consider first the simpler question of which minimal spanning set maximizes speed of learning.

Observe that if we restrict to any minimal spanning set (effectively, throwing out the remaining signals), then we return to a setting of exactly sufficient information. Let us generalize the previous analysis in the following way: for each minimal spanning set $\mathcal{S}$, define $\beta_i^{\mathcal{S}}$ to be the coefficients satisfying

$$e_1 = \sum_{i \in \mathcal{S}} \beta_i^{\mathcal{S}} \cdot c_i.$$

By Corollary 1, optimal signal acquisitions from any minimal spanning set $\mathcal{S}$ yields a posterior variance of approximately $\left( \sum_{i \in \mathcal{S}} |\beta_i^{\mathcal{S}}| \right)^2 / t$ at all large times $t$. Notice that posterior variance is strictly increasing in

$$\phi(\mathcal{S}) = \sum_{i \in \mathcal{S}} |\beta_i^{\mathcal{S}}|,$$

as it can be rewritten as $(\phi(\mathcal{S}))^2 / t$. Throughout, we work with the simpler statistic $\phi(\mathcal{S})$. The smaller $\phi(\mathcal{S})$ is, the faster the community learns.[14]

We maintain throughout the following assumption on the coefficient matrix $C$:

**Assumption 2** (Unique Minimizer). *$\phi(\mathcal{S})$ has a unique minimizer $\mathcal{S}^*$ among minimal spanning sets $\mathcal{S} \subset [N]$.*

This assumption says that there is a unique minimal spanning set that maximizes speed of learning. It fails in examples such as the following:

**Example 1.** *The signals are $X_1 = \omega + \epsilon_1$ and $X_2 = \omega + \epsilon_2$. Clearly these signals are duplicates of one another, and learning occurs equally fast from either of the minimal spanning sets $\{X_1\}$ or $\{X_2\}$.*

**Example 2.** *The signals are $X_1 = \omega + b_1 + \epsilon_1$, $X_2 = b_1 + \epsilon_2$, $X_3 = \omega + b_2 + \epsilon_3$, and $X_4 = b_2 + \epsilon_4$. In this environment, learning occurs equally fast from either of the minimal spanning sets $\{X_1, X_2\}$ and $\{X_3, X_4\}$.*

---

[14]We can extend this definition to arbitrary set of signals $\mathcal{A} \subset [N]$ (not necessarily minimally-spanning) as follows. For any set that contains a minimal spanning set, define

$$\phi(\mathcal{A}) = \min_{\mathcal{S} \subset \mathcal{A}} \phi(\mathcal{S}),$$

where the minimum is taken over all minimal spanning sets $\mathcal{S}$ contained in $\mathcal{A}$. If such $\mathcal{S}$ does not exist (i.e., $\mathcal{A}$ is not itself spanning), we let $\phi(\mathcal{A}) = \infty$. In particular, $\phi([N]) = \min_{\mathcal{S} \subset [N]} \phi(\mathcal{S})$ represents the minimum asymptotic standard deviation achieved by observing only those signals in some minimal spanning set.

These examples are special, and Assumption 2 holds when we are permitted arbitrarily small perturbations of the environments above.[15]

From here on, we work under this assumption, so that there is a "best" minimal spanning set $\mathcal{S}^*$. Define the frequency vector $\lambda^* \in \Delta^{N-1}$ by

$$\lambda_i^* = \begin{cases} \frac{|\beta_i^{\mathcal{S}^*}|}{\sum_{j \in \mathcal{S}^*} |\beta_j^{\mathcal{S}^*}|} & \forall \ i \in \mathcal{S}^* \\ 0 & \forall \ i \notin \mathcal{S}^* \end{cases} \tag{3}$$

Then, $\lambda^*$ is the optimal sampling rule when we restrict long run observations to sample exclusively from a single minimal spanning set. Our first theorem, now stated, says that $\lambda^*$ remains optimal even when we do not impose any restrictions on the sampling procedure; that is, so long as $C$ satisfies Unique Minimizer, then the best long run strategy is to restrict to the best minimal spanning set, and to sample from that set as in the previous section.

**Theorem 1.** *Suppose that the coefficient matrix $C$ satisfies Unique Minimizer, with $\mathcal{S}^*$ uniquely minimizing $\phi(\mathcal{S})$. Let $\lambda^*$ be given by (3). Then $n_i(t) \sim \lambda_i^* \cdot t$ for each signal $i$.*[16]

The conclusion can be loosely interpreted as stating that $\lambda^*$ is the "most efficient linear representation" of the payoff-relevant state in terms of the signal coefficients.[17]

The necessity of Assumption 2 for Theorem 1 is seen from the previous Example 1 and Example 2, which did not satisfy Unique Minimizer. In Example 1, all divisions across signals are (trivially) equally optimal. We show in Appendix A that it is possible for infinite observations of more than $K$ signals to be *strictly* optimal, using the environment given in Example 2.

Finally, we point out the following comparative static.

**Corollary 3.** *Suppose that the coefficient matrix $C$ satisfies Unique Minimizer. Write each signal as $X_i = \alpha \langle c_i, \theta \rangle + \epsilon_i$ , so that the precision of signal $X_i$ is increasing in $\alpha$. Then, either $\lambda_i^* = 0$ or $\lambda_i^*$ is locally decreasing in $\alpha$.*

---

[15] Assumption 2 is generically satisfied.

[16] We conjecture that the stronger conclusion $n_i(t) = \lambda_i^* \cdot t + O(1)$ also holds. In Remark 2 in the appendix, we prove this conjecture assuming $|\mathcal{S}^*| = K$.

[17] Specifically, consider the following constrained minimization problem: $\min \sum_{i=1}^N |\beta_i|$ subject to $\sum_{i=1}^N \beta_i \cdot c_i = e_1$. It can be shown by linear programming that the minimum is attained exactly when $\beta_i = \beta_i^*$—that is, when focusing on a single minimal spanning set.

That is, if signal $i$ is viewed with positive frequency in the long run, then its asymptotic frequency is (locally) decreasing in its precision. This implies loosely that a signal is most frequently viewed when it is *least informative* subject to being in the *most informative* set.

# 6 Main Results

We move on now to our main analysis: characterization of long-run acquisitions, and when these converge to the optimal acquisitions discussed above. In general, we may expect a difference between the best one-shot allocation of $t$ acquisitions, and the set of $t$ acquisitions that are chosen by sequential decision-makers. We show that whether society's acquisitions $m(t)$ eventually approximate the optimal acquisitions $n(t)$ depends critically on how many signals are required to identify $\omega$.

   We present below two versions of our main results. In Section 6.1, we restrict to a special class of environments, where the main results are simpler to state and the main intuition (for the general setting) is clearer to see. We then turn to the general setting in Section 6.2.

## 6.1 Special Class of Environments

Throughout this section, we impose the following assumption:

**Assumption 3** (Strong Linear Independence). *$N \geq K$, and every $K \times K$ submatrix of $C$ is of full rank.*

   Strong Linear Independence holds generically, but rules out environments such as the following:

**Example 3.** *Suppose that the available signals are $X_1 = \omega + b_1 + \epsilon_1$, $X_2 = b_1 + \epsilon_2$, $X_3 = \omega + b_2 + \epsilon_3$, $X_4 = b_2 + \epsilon_4$, $X_5 = \omega + b_3 + \epsilon_5$, and $X_6 = b_3 + \epsilon_6$. Then $K = 4$ and the signals $X_1, X_2, X_3, X_4$ are* not *linearly independent.*

   The example below shows that sequential information acquisition need not lead to the optimal sampling procedure; in particular, agents can become "stuck" observing signals from a sub-optimal spanning set.

**Example 4.** *There are three signals*

$$X_1 = \omega/2 + \epsilon_1$$
$$X_2 = \omega + b_1 + \epsilon_2$$
$$X_3 = \omega - b_1 + \epsilon_3$$

*Note that $\phi(\{X_1\}) = 2 > 1 = \phi(\{X_2, X_3\})$, so the latter two signals maximize speed of learning.*

*However, consider a prior belief such that $\omega$ and $b_1$ are independent, the variance about $\omega$ is 1, and the variance about $b_1$ exceeds 3. Prior to any observations, the precision of the first signal $\frac{\omega}{2} + b_1$ (interpreted as a noisy observation of $\omega$) is $\frac{1}{4}$, whereas the latter two signals $\omega + b_1 + \epsilon_2$ and $\omega - b_1 + \epsilon_3$ each has lower precision.[18] Thus the best choice in the first period is to observe $X_1$. Since this observation does not affect the variance of $b_1$, the same argument shows that* every *agent observes signal 1.*

The result below (stated as a corollary, since it will follow from Theorem 2 in the subsequent section) generalizes this example to show that different priors can lead to different absorbing sets.

**Corollary 4.** *Suppose Strong Linear Independence is satisfied. For any minimal-spanning set that contains fewer than $K$ signals, there exists an open set of prior beliefs under which agents exclusively observe signals from this set.*

We note that the implied inefficiency, measured as the ratio of the optimal speed of learning and the achieved speed of learning, can be an arbitrarily large constant. Specifically, for any positive number $L$, there exists an environment in which

$$\phi(\mathcal{S})/\phi(\mathcal{S}^*) > L$$

where $\mathcal{S}$ is the asymptotic observation set and $\mathcal{S}^*$ is the optimal asymptotic observation set. This can be shown by direct construction: modify the example above so that

$$X_1 = \omega/2 + \epsilon_1$$
$$X_2 = \alpha\omega + b_1 + \epsilon_2$$
$$X_3 = \alpha\omega - b_1 + \epsilon_3$$

---

[18]The signal $X_1 = \omega/2 + \epsilon_1$ is equivalent to the signal $\omega + 2\epsilon_1$, which has distribution $\mathcal{N}(\omega, 4)$. Each of $\omega + b_1 + \epsilon_2$ and $\omega + b_2 + \epsilon_2$ has greater variance conditional on $\omega$.

with $\alpha > \frac{L}{2}$. We note that the set of "inefficient" priors (which result in sub-optimal learning) does decrease in size as the level of inefficiency increases.[19]

Converse to the above result, our next result shows that if there is *no* minimal spanning set consisting of fewer than $K$ signals, then starting from *any* prior, information acquisition eventually concentrates on the best set of signals.

**Corollary 5.** *Suppose Strong Linear Independence and Unique Minimizer are satisfied. Then, if every minimal spanning set has size $K$, starting from any prior belief, it holds that $m_i(t) \sim \lambda_i^* \cdot t, \forall i$.*

One may argue that generically (for randomly drawn coefficient vectors), every minimal spanning set is of size $K$. However, this notion of genericity ignores the fact that in many economic situations, information sources are not determined by a random process. Indeed, if we expect that sources are endogenous to design or strategic motivations, then important informational environments may be "non-generic." For example, the existence of any source that directly reveals $\omega$ (that is, $X = \alpha\omega + \epsilon$) is non-generic in the probabilistic sense, but plausible in practice. Sets of signals that partition into different groups are also economically interesting but non-generic.[20] Our earlier Corollary 5 shows that inefficiency is a likely outcome in these cases.

The intuition for the above results, and in particular the role of "small" minimal spanning sets, is as follows. If a minimal spanning set consists of $K$ signals, then it must also have "full rank". As agents accumulate observations from any minimal spanning set, they learn not only about $\omega$ but also about all of the biases. The aggregated information in the community must then eventually swamp the prior, so that the asymptotic evaluation of the value of different signals at large periods $t$ is prior-independent. In fact, this asymptotic evaluation returns the optimal divisions in Section 5.

The argument above is no longer valid when there is a smaller set of signals that is minimally spanning. Intuitively, observation of $k < K$ sources can be self-reinforcing, since at most $k$ unknown states are revealed from these sources. Thus,

---

[19]As $\alpha$ increases, the prior variance about $b_1$ has to be sufficiently large in order for the first agent to choose $X_1$.

[20]We point out that the set of coefficient matrices that satisfy Unique Minimizer is "generic" in the following stronger sense: fixing the *directions* of coefficient vectors (as in Corollary 3), and suppose that the *precisions* are drawn at random, then different minimal spanning sets correspond to different speed of learning. In contrast, whether every minimal spanning set has size $K$ is a condition on the *directions* themselves, so this stronger notion of genericity does not apply.

any uncertainty in the prior about the other $K - k$ states can persist, despite infinitely many observations of the $k$ signals. If the remaining sources depend on these $K = k$ "poorly understood" states, then agents can become locked in observing the original $k$ sources.

Returning to our previous example of sequential choice of locations for RCTs, application of the results above yields two possible long run outcomes (depending on the pattern of correlation across signals). One outcome is that information acquisitions eventually approximate the best possible sampling rule over sites, and another is that information acquisitions eventually concentrate on sites that do not jointly maximize speed of learning. The key property that separates these two long-run outcomes is whether there is a set of locations that are *revealing*—so that given sufficiently many RCTs at those locations, the community will recover $\omega$—and moreover *self-contained*—so that they provide limited information that would help future development economists interpret estimates at other locations. If so, then long run learning need not be efficient. In contrast, if these conditions are not fulfilled, then optimal long run learning will obtain.

Although we defer the proof of Corollary 5 to the appendix, we conclude this section with a brief sketch of the formal argument. Our strategy is to work with the following function, defined on frequency vectors:

$$f^*(\lambda_1, \ldots, \lambda_N) = \lim_{t \to \infty} t \cdot f(\lambda_1 t, \ldots, \lambda_N t).$$

We show in an intermediate step that signal acquisitions chosen according to a frequency vector that minimizes $f^*$ will asymptotically also minimize the posterior variance function $f$. This justifies our analysis of $f^*$ (see Lemma 4).

The function $f^*$ is convex in $\lambda$, and its unique minimum turns out to be the vector $\lambda^*$ (see Lemma 6).[21] The question of whether optimal long-run acquisitions are achieved is equivalent to the question of whether signal acquisition frequencies converge to $\lambda^*$.

Society's acquisitions follow a procedure of "pseudo"-gradient descent, where the vector $\lambda(t) = m(t)/t$ evolves according to

$$\lambda(t+1) = \frac{t}{t+1} \lambda(t) + \frac{1}{t+1} e_i.$$

The vector $e_i$ is the coordinate vector that yields the greatest (immediate) reduction in $f$ (and roughly the greatest reduction in $f^*$). Unlike standard gradient descent,

---

[21]Here we use the assumption that $\mathcal{S}^*$ is the uniquely "best" minimal spanning set.

the descent here can occur only along a finite set of feasible directions, corresponding to the available signals.

However, this limitation is without loss whenever $f^*$ is differentiable at $\lambda$, since then all directional derivatives can be rewritten as a convex combination of partial derivatives along basis vectors. We observe that $f^*$ is differentiable whenever $\lambda$ has at least $K$ nonzero coordinates.

It remains to show that $\lambda(t)$ will in fact eventually have at least $K$ nonzero coordinates (corresponding to positive frequency of observation of at least $K$ signals). To show this, we argue that society must eventually observe *some* minimal spanning set, or else it would not recover $\omega$. Combining this with the above, if every minimal spanning set has size $K$, then descent is well-behaved and ends at the global minimum $\lambda^*$. This yields Corollary 5, and Theorem 3 (in the subsequent section) follows from a similar argument.

In contrast, pseudo-gradient descent can become "stuck" at vectors $\lambda$ with fewer than $K$ nonzero coordinates. Formally, $f^*$ can fail to be differentiable at these points. If that is the case, observation of $k < K$ sources can be self-reinforcing, as reflected in Corollary 4 and more generally Theorem 2 below.

## 6.2 General Case

Towards the result for the general setting, we introduce the notion of *subspace optimality*. For any spanning set of signals $\mathcal{S}$, let $\overline{\mathcal{S}} \subseteq [N]$ be the set of available signals whose coefficient vectors belong to the subspace spanned by signals in $\mathcal{S}$. Notice in particular that $\overline{\mathcal{S}}$ contains $\mathcal{S}$. We say a minimal spanning set $\mathcal{S}$ is *subspace-optimal* if it uniquely minimizes $\phi$ among available subsets of $\overline{\mathcal{S}}$. For example, if the available signals are $X_1 = \omega/2 + \epsilon_1$ and $X_2 = \omega + \epsilon_2$, then $\{X_1\}$ is a minimal spanning set, but it is not optimal in its subspace.[22]

**Theorem 2.** *Suppose $\mathcal{S}$ is a minimal spanning set that is moreover subspace-optimal. Then, there exists an open set of prior beliefs under which long-run frequencies are strictly positive for signals in $\mathcal{S}$, and zero everywhere else.*[23]

Notice that under the earlier assumption of Strong Linear Independence, any minimal spanning set with fewer than $K$ signals is subspace-optimal.[24] Thus, as we

---

[22] $X_2$ belongs to the subspace spanned by $X_1$, and $\phi(\{X_2\}) < \phi(\{X_1\})$.

[23] These frequencies are the optimal frequencies when only signals in $\mathcal{S}$ are available (i.e., given by Corollary 1).

[24] Any other signal that belongs to its subspace would violate linear independence.

saw in Corollary 4, the possibility of inefficiency hinged on existence of such sets.

Converse to Theorem 2, our next result shows that starting from *any* prior, information acquisition eventually concentrates on a set of signals that is subspace-optimal. We use an assumption which strengthens Unique Minimizer.

**Assumption 4** (Unique Minimizer in Every Subspace). *For any spanning set $\mathcal{A} \subset [N]$, $\operatorname{argmin}_{\mathcal{S} \subset \overline{\mathcal{A}}} \phi(\mathcal{S})$ has a unique solution, where the minimum is taken over minimal spanning sets $\mathcal{S}$.*

This says that in every spanning subspace, there exists a unique minimal spanning subset $\mathcal{S}$ that maximizes asymptotic speed of learning. It is clearly guaranteed if different minimal spanning sets correspond to different values of $\phi$.

**Theorem 3.** *Suppose that Assumption 4 is satisfied. Given any prior belief, long-run frequencies exist for every signal. Moreover, if $\mathcal{S}$ denotes the signals viewed with positive frequencies, then $\mathcal{S}$ is a minimal spanning set that is subspace-optimal.*

Notice that if every minimal spanning set is of size $K$, then all minimal spanning sets belong to the same subspace. Furthermore, if Unique Minimizer holds, there can only be one minimal spanning set that is optimal in its subspace, and moreover this is the "best" set (in the sense of Section 5). This yields the previous Corollary 5 from the theorem above.

# 7   Information Interventions

Section 6 demonstrated the possibility for sequential information acquisition to result in inefficient learning. We ask now whether it is possible for a benevolent outside party to help society achieve efficient learning by providing a one-time injection of free information. Naturally, this question applies only when agents (on their own) could eventually achieve a sub-optimal speed of learning. The conditions under which this occurs are given in Theorem 2.

Formally, suppose a policy-maker chooses $M$ signals $\langle p_j, \theta \rangle + \mathcal{N}(0, 1)$, where each $\|p_j\|_2 \leq \gamma$, so that signal precisions are bounded by $\gamma^2$. At time 0, this information is made public. All subsequent agents update their prior beliefs based on this free information, and also on the history of signal acquisitions thus far. The goal of the policy-maker is to maximize the community's asymptotic speed of learning. Below, we use *efficient learning* to mean the case in which the asymptotic speed of learning

achieves the optimum—that is, the final observation set is $\mathcal{S}^*$ and long-run frequencies are $\lambda^*$.

Is there a sufficient number of (kinds of) signals, such that efficient learning can be guaranteed? We answer in the affirmative below: $K-1$ precise signals are sufficient to produce efficient learning:

**Proposition 2.** *For any prior, there exist $\gamma$ and $K-1$ signals with $\|p_j\|_2 \leq \gamma$ such that with these free signals provided at $t = 0$, society achieves efficient learning.*

Intuitively, as long as the free signals make agents sufficiently informed about the biases $b_1, \ldots, b_{K-1}$, they can preclude the situation in which agents get stuck in a sub-optimal set as in Example 4. Notice that optimal information intervention does not need to teach directly about $\omega$ (the payoff-relevant state), which the agents will learn on their own. Rather, the planner should only provide auxiliary information that helps agents to better interpret the sources.

# 8 Related Literature

In addition to the references mentioned in the introduction, our results build on prior work regarding speed of learning (Vives, 1992; Golub and Jackson, 2012; Harel et al., 2017; Hann-Caruthers, Martynov and Tamuz, 2017), and is related also to the experimental design literature in statistics (see Chernoff (1972) for a survey). Specifically, our results in Section 5 are related to **c**-optimality, in which $t$ experiments are chosen to minimize the posterior variance of a linear combination of the unknown states (in our case, simply the posterior variance of the first unknown state). Theorem 1 can be seen as an integer design version of the problem considered in Chaloner (1984). Chaloner (1984) showed that a $c$-optimal Bayesian continuous design exists on at most $K$ points, but does not provide a construction of this design. Extending this, we supply a characterization of the optimal design itself; this improves on the prior result by showing uniqueness of the optimal design, and demonstrating that for certain correlational structures, optimal design exists on strictly fewer than $K$ points.

# 9 Conclusion

We study a model of sequential learning, where agents choose what kind of information to acquire from a large set of information sources. The key force of interest is

the externality that current informational choices generate on future agents.

Our main results characterize two starkly different possibilities and the conditions under which either obtains: (1) the externality is *beneficial*: past information acquisitions help future agents to discern which sources are most informative, and in the long run, agents converge to acquiring information only from the most informative sources; (2) the externality is *harmful*: past information acquisitions increase the value of "low-quality" sources relative to "high-quality" sources, pushing future agents to acquire information from a set of sources that yields inefficiently slow learning. A simple property of the correlation structure across sources determines when such "learning traps" emerge, and which sources are a part of them.

When a community is stuck observing inefficient sources, what kind of information interventions might push the community towards efficient learning? One possibility is to limit the number of sources, and especially to remove "decoy" sources that are low-quality but self-reinforcing. Another possibility is to provide agents with free information. We show that a policy-maker can guarantee efficient long-run learning if he provides a sufficient number of sufficiently precise signals. The optimal information intervention does not inform directly about the payoff-relevant state, but rather provides auxiliary information that helps agents to interpret the best sources (so that these are subsequently observed). This intervention may require educating agents along *many* different dimensions: we conjecture that provision of a single kind of information (no matter how precise) can be ineffective in a large number of environments. This points to the potential long-run ineffectiveness of information campaigns that are very informative but limited in scope.

Finally, although in this paper we focus on informational demand given a *fixed* set of information sources, one may also consider the reverse question of what kinds of information will be *endogenously* provided by strategic sources. Our results suggest that the answer to this question can be subtle: information sources most frequently viewed in the long run are those that are "least informative in a most informative set." Thus, a source that wants to maximize frequency of viewership has two competing incentives: first, to be viewed at all within the competitive market, it must provide sufficiently useful information; second, conditional on being viewed, it wants to reveal information slowly (so as to increase the number of observations). We leave characterization of the supply of information in an "informationally overabundant" environment for future work.

# References

**Ali, Nageeb.** 2017. "Herding with Costly Information." Working Paper.

**Banerjee, Abhijit.** 1992. "A Simple Model of Herd Behavior." *Quaterly Journal of Economics*, 107(3): 797–817.

**Bikhchandani, Sushil, David Hirshleifer, and Ivo Welch.** 1992. "A Theory of Fads, Fashion, Custom, and Cultural Change as Information Cascades." *Journal of Political Economy*, 100(5): 992–1026.

**Burguet, Roberto, and Xavier Vives.** 2000. "Social Learning and Costly Information." *Economic Theory.*

**Chaloner, Kathryn.** 1984. "Optimal Bayesian Experimental Design for Linear Models." *The Annals of Statistics*, 12(1): 283–300.

**Chernoff, Herman.** 1972. *Sequential Analysis and Optimal Design.* Society for Industrial and Applied Mathematics.

**Che, Yeon-Koo, and Konrad Mierendorff.** 2017. "Optimal Sequential Decision with Limited Attention." Working Paper.

**Fudenberg, Drew, Philip Strack, and Tomasz Strzalecki.** 2017. "Stochastic Choice and Optimal Sequential Sampling." Working Paper.

**Golub, Benjamin, and Matthew Jackson.** 2012. "How Homophily Affects the Speed of Learning and Best-Response Dynamics." *The Quarterly Journal of Economics.*

**Hann-Caruthers, Wade, Vadim Martynov, and Omer Tamuz.** 2017. "The Speed of Sequential Asymptotic Learning." Working Paper.

**Harel, Matan, Elchanan Mossel, Philipp Strack, and Omer Tamuz.** 2017. "The Speed of Social Learning." Working Paper.

**Liang, Annie, Xiaosheng Mu, and Vasilis Syrgkanis.** 2017. "Dynamic Information Acquisition from Multiple Sources." Working Paper.

**Mayskaya, Tatiana.** 2017. "Dynamic Choice of Information Sources." Working Paper.

**Mueller-Frank, Manuel, and Mallesh Pai.** 2016. "Social Learning with Costly Search." *American Economic Journal: Microeconomics.*

**Sethi, Rajiv, and Muhamet Yildiz.** 2016. "Communication with Unknown Perspectives." *Econometrica*, 84(6): 2029–2069.

**Sethi, Rajiv, and Muhamet Yildiz.** 2017. "Culture and Communication." Working Paper.

**Smith, Lones, and Peter Sorenson.** 2000. "Pathological Outcomes of Observational Learning." *Econometrica.*

**Vives, Xavier.** 1992. "How Fast do Rational Agents Learn?" *Review of Economic Studies.*

# A  Examples Failing Unique Minimizer

## A.1  First Example

**Example 5.** *There are $K = 3$ states $\omega, b_1, b_2$ independently drawn with prior variances $\frac{1}{\alpha}, \frac{1}{\beta}, \frac{1}{\gamma}$. $N = 4$ signals are available, and they are respectively*

$$X_1 = \omega + b_1 + \epsilon_1$$
$$X_2 = b_1 + \epsilon_2$$
$$X_3 = \omega + b_2 + \epsilon_3$$
$$X_4 = b_2 + \epsilon_4$$

*with standard normal errors. Note that the former two signals and the latter two signals are both spanning, and these two sets generate the same asymptotic variance. Thus Assumption 2 is not satisfied.*

*The posterior variance about $\omega$ as a function of the number of observations $q_1, q_2, q_3, q_4$ of each signal type can be derived as follows. First, given $q_2$ observations of signal $X_2$ and $q_4$ observations of signal $X_4$, posterior variance about $\theta_2$ and $\theta_3$ are $1/(q_2 + \beta)$ and $1/(q_4 + \gamma)$ respectively. Consider now $q_1$ additional observations of $X_1$; this provides the same information about the payoff-relevant state $\omega$ as the signal $\omega + \epsilon'$, where $\epsilon'$ is an independent noise term with variance $\frac{1}{q_1} + \frac{1}{q_2+\beta}$. Similarly, $q_3$ additional observations of $X_3$ are equivalent to the signal $\omega + \epsilon''$, where $\epsilon''$ is an independent noise term with variance $\frac{1}{q_3} + \frac{1}{q_4+\gamma}$. From this we deduce that posterior variance about $\omega$ is*

$$f(q_1, q_2, q_3, q_4) = 1 \left/ \left( \alpha + \frac{1}{\frac{1}{q_1} + \frac{1}{q_2+\beta}} + \frac{1}{\frac{1}{q_3} + \frac{1}{q_4+\gamma}} \right) \right. .$$

*The optimal division vector thus seeks to* maximize

$$\frac{1}{\frac{1}{q_1} + \frac{1}{q_2+\beta}} + \frac{1}{\frac{1}{q_3} + \frac{1}{q_4+\gamma}} \tag{4}$$

*It is useful to rewrite (4) in the following way:*

$$\frac{1}{4} \left( q_1 + q_2 + \beta + q_3 + q_4 + \gamma - \frac{(q_1 - q_2 - \beta)^2}{q_1 + q_2 + \beta} - \frac{(q_3 - q_4 - \gamma)^2}{q_3 + q_4 + \gamma} \right) .$$

*Then, since $q_1 + q_2 + \beta + q_3 + q_4 + \gamma = t + \beta + \gamma$ is fixed at any time $t$, it is equivalent to choose $q_1, q_2, q_3, q_4$ to minimize the sum of ratios*

$$\frac{(q_1 - q_2 - \beta)^2}{q_1 + q_2 + \beta} + \frac{(q_3 - q_4 - \gamma)^2}{q_3 + q_4 + \gamma} .$$

*Ideally, if signals were perfectly divisible, the optimum would be to choose $q_1 = q_2 + \beta$ and $q_3 = q_4 + \gamma$. But as each $q_i$ is restricted to integer values, this continuous optimum is not feasible whenever $\beta$ and $\gamma$ are not integers.*

*The solution to this integer optimization problem is involved, and the details are relegated to Appendix A.2. To express the solution, we need some additional notation. Let $r$ be the integer that minimizes $|r - \beta|$ (the distance to $\beta$) and let $s$ be the integer that minimizes $|s - \gamma|$. Further, let $\langle\beta\rangle$ and $\langle\gamma\rangle$ be the value of these absolute differences. We show that when the parity of $t$ and $r + s$ are the same, the optimal $(q_1, q_2, q_3, q_4)$ satisfy*

$$q_1, q_2 \approx \frac{\langle\beta\rangle}{2\langle\beta\rangle + 2\langle\gamma\rangle} \cdot t; \quad q_3, q_4 \approx \frac{\langle\gamma\rangle}{2\langle\beta\rangle + 2\langle\gamma\rangle} \cdot t.$$

*and otherwise the optimal $(q_1, q_2, q_3, q_4)$ satisfy*

$$q_1, q_2 \approx \frac{\langle\beta\rangle}{2\langle\beta\rangle + 2 - 2\langle\gamma\rangle} \cdot t; \quad q_3, q_4 \approx \frac{1 - \langle\gamma\rangle}{2\langle\beta\rangle + 2 - 2\langle\gamma\rangle} \cdot t.$$

*Thus, all four signals are observed with positive frequencies in the long run according to the optimal criterion.*

Although the example is involved, its intuition is simple: we would most like to set $q_1 = q_2 + \beta$ and $q_3 = q_4 + \gamma$, but this is not feasible when $\beta$ and $\gamma$ are not integers. Thus, there is inevitably some loss from the ideal case where signals are continuously divisible. This loss turns out to be convex in signal counts, so to minimize total loss, both groups of signals are observed infinitely often.

The conclusion of Theorem 1 fails to hold in a strong sense in the example above, since *all* signals are observed infinitely often. Later we provide another example that does not satisfy Unique Minimizer, but where the conclusion of Theorem 1 holds "qualitatively." The difference in these two examples, and in addition the complexity of derivation of the asymptotic frequencies above suggest that characterization of optimal acquisitions is in general difficult without the Unique Minimizer assumption.

## A.2 Details for Example 5

To solve the integer maximization problem (4), let $r$ be the integer that minimizes $|r - \beta|$ (the distance to $\beta$) and let $s$ be the integer that minimizes $|s - \gamma|$. Further, let $\langle\beta\rangle$ and $\langle\gamma\rangle$ be the value of these absolute differences. We assume $2\beta, 2\gamma$ are not integers, so that $0 < \langle\beta\rangle, \langle\gamma\rangle < \frac{1}{2}$. We also assume $\langle\beta\rangle \neq \langle\gamma\rangle$, and by symmetry focus on the case of $\langle\beta\rangle < \langle\gamma\rangle$.

With these assumptions, it is clear that when $q_1, q_2$ are integers, the minimum value of $|q_1 - q_2 - \beta|$ is $\langle\beta\rangle$, achieved if and only if $q_1 = q_2 + r$. Similarly the minimum value of $|q_3 - q_4 - \gamma|$ is $\langle\gamma\rangle$, achieved when $q_3 = q_4 + s$. Now if the total number of observations $t$ has the *same parity* as $r + s$, it is possible to choose $q_1, q_2, q_3, q_4$ such that their sum is $t$ and $q_1 = q_2 + r$, $q_3 = q_4 + s$—any pair $q_2, q_4$ with sum $\frac{t-r-s}{2}$ leads to such a solution. Given these constraints, then, the optimum is to choose $q_2, q_4$ to minimize $\frac{\langle\beta\rangle^2}{2q_2+r+\beta} + \frac{\langle\gamma\rangle^2}{2q_4+s+\gamma}$. The optimal $q_2$ and $q_4$ satisfy $q_2/q_4 \approx \langle\beta\rangle/\langle\gamma\rangle$, which together with $q_2 + q_4 = \frac{t-r-s}{2}$ implies

$$q_1, q_2 \approx \frac{\langle\beta\rangle}{2\langle\beta\rangle + 2\langle\gamma\rangle} \cdot t; \quad q_3, q_4 \approx \frac{\langle\gamma\rangle}{2\langle\beta\rangle + 2\langle\gamma\rangle} \cdot t.$$

27

On the other hand, suppose $t$ has the *opposite parity* to $r+s$. In this case $q_1 = q_2 + r$ and $q_3 = q_4 + s$ cannot both hold, thus $|q_1 - q_2 - \beta|$ and $|q_3 - q_4 - \gamma|$ cannot both take their minimum values $\langle\beta\rangle$ and $\langle\gamma\rangle$. It turns out that the best one can do is choose $q_1 = q_2 + r$ and $q_3 = q_4 + s \pm 1$ so that $|q_1 - q_2 - \beta| = \langle\beta\rangle$ and $|q_3 - q_4 - \gamma| = 1 - \langle\gamma\rangle$. Then, the optimal choice of $q_2, q_4$ with sum $\frac{t-r-s\mp 1}{2}$ to minimize $\frac{\langle\beta\rangle^2}{2q_2+r+\beta} + \frac{(1-\langle\gamma\rangle)^2}{2q_4+s+\gamma\pm 1}$. This yields

$$q_1, q_2 \approx \frac{\langle\beta\rangle}{2\langle\beta\rangle + 2 - 2\langle\gamma\rangle} \cdot t; \quad q_3, q_4 \approx \frac{1 - \langle\gamma\rangle}{2\langle\beta\rangle + 2 - 2\langle\gamma\rangle} \cdot t.$$

## A.3 Second Example

In the following example, Unique Minimizer is violated. However, the qualitative conclusion of Theorem 1 still holds. Namely, as $t \to \infty$, at most $K$ signals are observed with positive frequency under the $t$-optimal division.

**Example 6.** *Consider state $\omega$ and bias $b_1$ (prior beliefs will be specified shortly). There are three signals $\omega + b_1 + \epsilon_1$, $\omega - b_1 + \epsilon_2$ and $\omega + \epsilon_3$, where each noise term is standard normal. We assume the prior beliefs are such that $\omega + b_1$ and $\omega - b_1$ are independent, with variances $\frac{1}{\alpha}$ and $\frac{1}{\beta}$. Observe that $\phi(\{1,2\}) = 1 = \phi(\{3\})$, so Unique Minimizer fails.*

*We claim that whenever $\alpha - \beta$ is not an integer, $t$-optimal divisions choose the third signal only a bounded number of times. Intuitively, this is because one observation of $\omega + b_1 + \epsilon_1$ combined with one observation of $\omega - b_1 + \epsilon_2$ contain at least as much information as their sum $2\omega + \epsilon_1 + \epsilon_2$, which is equivalent to two observations of $\omega + \epsilon_3$. Thus, devoting any level of attention to the third signal is* weakly *worse than splitting that attention evenly between the first two signals. Moreover, the combination of the first two signals also informs about $b_1$, which is correlated with the payoff-relevant state $\omega$ whenever $\alpha \neq \beta$. Thus, society optimally "ignores" the third signal if its (prior and posterior) beliefs about $\omega + b_1$ and $\omega - b_1$ are asymmetric. As we show below, this occurs precisely when $\alpha - \beta$ is not an integer.*

*To formalize the above intuition, we observe that given $q_1$ observations of signal 1 and $q_2$ observations of signal 2, society's posterior variance about $\omega$ is $\left(\frac{1}{q_1+\alpha} + \frac{1}{q_2+\beta}\right)/4$. Thus, with $q_3$ additional observations of the third signal, society's posterior variance becomes*

$$f(q_1, q_2, q_3) = 1 \left/ \left( \frac{4}{\frac{1}{q_1+\alpha} + \frac{1}{q_2+\beta}} + q_3 \right) \right. .$$

*The optimal division at time $t$ thus maximizes*

$$\max_{q_1,q_2,q_3 \in \mathbb{Z}^+, q_1+q_2+q_3=t} \frac{4}{\frac{1}{q_1+\alpha} + \frac{1}{q_2+\beta}} + q_3.$$

*The maximand can be rewritten as*

$$\frac{4}{\frac{1}{q_1+\alpha} + \frac{1}{q_2+\beta}} + q_3 = q_1 + \alpha + q_2 + \beta + q_3 - \frac{(q_1 + \alpha - q_2 - \beta)^2}{q_1 + \alpha + q_2 + \beta}.$$

*Note that $q_1 + \alpha + q_2 + \beta + q_3 = t + \alpha + \beta$ is fixed, so society chooses $q_1, q_2$ to minimize the ratio $\frac{(q_1+\alpha-q_2-\beta)^2}{q_1+\alpha+q_2+\beta}$.*

*Suppose $\alpha - \beta$ is not an integer, let $\langle \alpha - \beta \rangle$ denote its distance to the nearest integer. Then, as $q_1, q_2$ are restricted to integers, the difference $|q_1 + \alpha - q_2 - \beta|$ takes minimum value $\langle \alpha - \beta \rangle > 0$. It follows that $\frac{(q_1 + \alpha - q_2 - \beta)^2}{q_1 + a + q_2 + b}$ is uniquely minimized by choosing $q_1, q_2$ such that $|q_1 + \alpha - q_2 - \beta| = \langle \alpha - \beta \rangle$ and $q_1 + q_2$ is as large as possible. Hence, both $q_1$ and $q_2$ are close to $\frac{t}{2}$. As we claimed, t-optimal division eventually focuses on the first two signals.*

# B   Posterior Variance Function

## B.1   A Basic Lemma

Here we review and extend a basic result from Liang, Mu and Syrgkanis (2017). Specifically, we show that the posterior variance about $\omega$ weakly decreases over time, and the marginal value of any signal decreases in its signal count.

**Lemma 2.** *Given prior covariance matrix $V^0$ and $q_i$ observations of each signal $i$, society's posterior variance about $\omega$ is given by*

$$f(q_1, \ldots, q_N) = \left[ ((V^0)^{-1} + C'QC)^{-1} \right]_{11} \tag{5}$$

*where $Q = \mathrm{diag}(q_1, \ldots, q_N)$. The function $f$ is decreasing and convex in each $q_i$ whenever these arguments take non-negative real values.*

*Proof.* Note that $(V^0)^{-1}$ is the prior precision matrix, and $C'QC = \sum_{i=1}^{N} q_i \cdot [c_i c_i']$ is the total precision from the signals. Thus (5) simply represents the fact that for Gaussian prior and signals, the posterior precision matrix is the sum of prior and signal precision matrices. To prove the monotonicity of $f$, consider the partial order $\succeq$ on positive semi-definite matrices where $A \succeq B$ if and only if $A - B$ is positive semi-definite. As $q_i$ increases, the matrix $Q$ and $C'QC$ increase in this order. Thus the posterior covariance matrix $((V^0)^{-1} + C'QC)^{-1}$ decreases in this order, which implies that the posterior variance about $\omega$ decreases. Intuitively, more information always improves the decision-maker's estimates.

To prove $f$ is convex, it suffices to prove $f$ is *midpoint-convex* since the function is clearly continuous. Take $q_1, \ldots, q_N, r_1, \ldots, r_N \in \mathbb{R}_+$ and let $s_i = \frac{q_i + r_i}{2}$. Define the corresponding diagonal matrices to be $Q$, $R$, $S$. Observe that $Q + R = 2S$. Thus by the AM-HM inequality for positive-definite matrices, we have in matrix order

$$((V^0)^{-1} + C'QC)^{-1} + ((V^0)^{-1} + C'RC)^{-1} \succeq 2((V^0)^{-1} + C'SC)^{-1}.$$

Using (5), we conclude

$$f(q_1, \ldots, q_N) + f(r_1, \ldots, r_N) \geq 2f(s_1, \ldots, s_N).$$

This proves the convexity of $f$. $\qquad\square$

## B.2 Inverse of Positive Semi-definite Matrices

For future use, we provide a definition of $[X^{-1}]_{11}$ for positive *semi-definite* matrices $X$. When $X$ is positive definite, its eigenvalues are strictly positive, and its inverse matrix is defined as usual. In general, we can apply the spectrum theorem to write

$$X = UDU'$$

with $U$ being a $K \times K$ orthogonal matrix whose columns are eigenvectors of $X$, and $D$ being a $K \times K$ diagonal matrix consisting of non-negative eigenvalues. Even if some of these eigenvalues are zero, we can think of $X^{-1}$ as

$$X^{-1} = (UDU')^{-1} = UD^{-1}U' = \sum_{j=1}^{K} \frac{1}{d_j} \cdot [u_j u_j']$$

with $u_j$ being the $j$-th column vector of $U$. We thus define

$$[X^{-1}]_{11} = \sum_{j=1}^{K} \frac{(\langle u_j, e_1 \rangle)^2}{d_j}, \tag{6}$$

with the convention that $\frac{0}{0} = 0$. Note that by this definition,

$$[X^{-1}]_{11} = \lim_{\epsilon \to 0_+} \left( \sum_{j=1}^{K} \frac{(\langle u_j, e_1 \rangle)^2}{d_j + \epsilon} \right) = [(X + \epsilon I_K)^{-1}]_{11}$$

since the matrix $X + \epsilon I_K$ has the same set of eigenvectors as $X$, with eigenvalues increased by $\epsilon$. Hence our definition of $[X^{-1}]_{11}$ is a continuous extension of the usual definition to positive semi-definite matrices. Note that we allow $[X^{-1}]_{11}$ to be infinite.

# C   Proof of Theorem 1

## C.1 Asymptotic Behavior of Posterior Variance

We first approximate the posterior variance as a function of the frequencies with which each signal is observed. Specifically,

**Lemma 3.** *For any $\lambda_1, \ldots, \lambda_N \geq 0$, let $\Lambda = \text{diag}(\lambda_1, \ldots, \lambda_N)$. Then*

$$f^*(\lambda_1, \ldots, \lambda_N) := \lim_{t \to \infty} t \cdot f(\lambda_1 t, \ldots, \lambda_N t)$$
$$= [(C'\Lambda C)^{-1}]_{11} \tag{7}$$

*Note that the matrix $C'\Lambda C$ is positive semi-definite. So the value of $[(C'\Lambda C)^{-1}]_{11}$ is well defined, see (6).*

*Proof.* Recall that $f(q_1, \ldots, q_N) = \left[ ((V^0)^{-1} + C'QC)^{-1} \right]_{11}$ with $Q = \text{diag}(q_1, \ldots, q_N)$. Thus

$$t f(\lambda_1 t, \ldots, \lambda_N t) = \left[ \left( \frac{1}{t}(V^0)^{-1} + C'\Lambda C \right)^{-1} \right]_{11}.$$

Hence by the continuity of $[X^{-1}]_{11}$ in the matrix $X$, we obtain the lemma. $\square$

We note that $C'\Lambda C$ is the Fisher Information Matrix when the signals are observed according to frequencies $\lambda$. Thus the above lemma can also be seen as an application of the Bayesian Central Limit Theorem.

## C.2 Reduction to the Study of $f^*$

The development of the function $f^*$ is useful for the following reason:

**Lemma 4.** *Suppose $\hat{\lambda}$ uniquely minimizes $f^*(\lambda)$ subject to $\lambda \in \Delta^{N-1}$ (the $N-1$-dimensional simplex), then the $t$-optimal divisions satisfy $n_i(t) \sim \hat{\lambda}_i \cdot t$ for each $i$.*

*Proof.* Fix any increasing sequence of times $t_1, t_2, \ldots$. It suffices to show that whenever the limit $\lambda_i := \lim_{m \to \infty} \frac{n_i(t_m)}{t_m}$ exists for each $i$, this limit $\lambda$ must be $\hat{\lambda}$. Suppose not, then by assumption $f^*(\lambda) > f^*(\hat{\lambda})$. For $\epsilon > 0$, define another vector $\tilde{\lambda} \in \mathbb{R}_+^N$ with $\tilde{\lambda}_i = \lambda_i + \epsilon, \forall i$. By the continuity of $f^*$, it holds that $f^*(\tilde{\lambda}) > f^*(\hat{\lambda})$ for sufficiently small $\epsilon$.

Since $\lambda_i = \lim_{m \to \infty} \frac{n_i(t_m)}{t_m}$, there exists $M$ sufficiently large such that $n_i(t_m) \leq \tilde{\lambda}_i \cdot t_m$ for each $i$ and $m \geq M$. Hence, for $m \geq M$,

$$t_m \cdot f(n_1(t_m), \ldots, n_N(t_m)) \geq t_m \cdot f(\tilde{\lambda}_1 \cdot t_m, \ldots, \tilde{\lambda}_N \cdot t_m) \to f^*(\tilde{\lambda}_1, \ldots, \tilde{\lambda}_N)$$

The first inequality uses the monotonicity of $f$. On the other hand,

$$t_m \cdot f(\hat{\lambda}_1 \cdot t_m, \ldots, \hat{\lambda}_N \cdot t_m) \to f^*(\hat{\lambda}_1, \ldots, \hat{\lambda}_N).$$

Comparing the above two displays, we see that for sufficiently large $m$, $f(n_1(t_m), \ldots, n_K(t_m)) > f(\hat{\lambda}_1 \cdot t_m, \ldots, \hat{\lambda}_N \cdot t_m)$. But this contradicts the $t$-optimality of the division $n(t_m)$, as society could do better by following frequencies $\hat{\lambda}$. The lemma is thus proved. $\square$

## C.3 Crucial Lemma

We pause to demonstrate the following technical lemma:

**Lemma 5.** *Suppose $\mathcal{S}^* = \{1, \ldots, K\}$ uniquely minimizes $\phi(\mathcal{S})$ and let $C^*$ be the $K \times K$ submatrix of $C$ corresponding to the first $K$ signals. Further suppose $\beta_j^{\mathcal{S}^*} = [(C^*)^{-1}]_{1j}$ is positive for $1 \leq j \leq K$. Then for any signal $i > K$, we can write $c_i = \sum_{j=1}^{K} \alpha_j \cdot c_j$ with $|\sum_{j=1}^{K} \alpha_j| < 1$.*

*Proof.* By assumption, we have the vector identity

$$e_1 = \sum_{j=1}^{K} \beta_j \cdot c_j \quad \text{with } \beta_j = [(C^*)^{-1}]_{1j} > 0.$$

Suppose for contradiction that $\sum_{j=1}^{K} \alpha_j \geq 1$ (the opposite case where the sum is $\leq -1$ can be similarly treated). In particular, some $\alpha_j$ is positive. Without loss of generality, we assume $\frac{\alpha_1}{\beta_1}$ is the largest among such ratios. Then $\alpha_1 > 0$ and

$$e_1 = \sum_{j=1}^{K} \beta_j \cdot c_j = \left( \sum_{j=2}^{K} (\beta_j - \frac{\beta_1}{\alpha_1} \cdot \alpha_j) \cdot c_j \right) + \frac{\beta_1}{\alpha_1} \cdot \left( \sum_{j=1}^{K} \alpha_j \cdot c_j \right)$$

This represents $e_1$ as a linear combination of the vectors $c_2, \ldots, c_K$ and $c_i$, with coefficients $\beta_2 - \frac{\beta_1}{\alpha_1} \cdot \alpha_2, \ldots, \beta_K - \frac{\beta_1}{\alpha_1} \cdot \alpha_K$ and $\frac{\beta_1}{\alpha_1}$. Observe that these coefficients are non-negative: for each $2 \leq j \leq K$, $\beta_j - \frac{\beta_1}{\alpha_1} \cdot \alpha_j$ is clearly positive if $\alpha_j \leq 0$ (since $\beta_j > 0$). And if $\alpha_j > 0$, then by assumption $\frac{\alpha_j}{\beta_j} \leq \frac{\alpha_1}{\beta_1}$ and $\beta_j - \frac{\beta_1}{\alpha_1} \cdot \alpha_j$ is again non-negative.

By definition, $\phi(\{2, \ldots, K, i\})$ is the sum of the absolute value of these coefficients. This sum is

$$\sum_{j=2}^{K} (\beta_j - \frac{\beta_1}{\alpha_1} \cdot \alpha_j) + \frac{\beta_1}{\alpha_1} = \sum_{j=1}^{K} \beta_j + \frac{\beta_1}{\alpha_1} \cdot (1 - \sum_{j=1}^{K} \alpha_j) \leq \sum_{j=1}^{K} \beta_j.$$

But then $\phi(\{2, \ldots, K, i\}) \leq \phi(\{1, 2, \ldots, K\})$, leading to a contradiction. Hence the lemma must be true. $\qquad\square$

## C.4 Proof of Theorem 1 when $|\mathcal{S}^*| = K$

Given Lemma 4, Theorem 1 will follow once we show that $\lambda^*$ uniquely minimizes $f^*(\lambda)$ over the simplex—recall that $\lambda^*$ denotes the optimal asymptotic frequencies for the minimal spanning set $\mathcal{S}^*$ that minimizes $\phi$. In this section, we prove $\lambda^*$ is indeed the unique minimizer whenever this "best" subset $\mathcal{S}^*$ contains exactly $K$ signals. Later on we will prove the same result even when $|\mathcal{S}^*| < K$, but that proof will require additional techniques.

**Lemma 6.** *Suppose $\mathcal{S}^* = \{1, \ldots, K\}$ is the unique minimizer of $\phi(\mathcal{S})$ over minimal spanning sets. Define $\lambda^* \in \Delta^{N-1}$ by*

$$\lambda_i^* = \frac{|[(C^*)^{-1}]_{1i}|}{\sum_{j=1}^{K} |[(C^*)^{-1}]_{1j}|}, 1 \leq i \leq K$$

*with $C^* = C_{[K][K]}$,[25] and $\lambda_i^* = 0, \forall i > K$. Then $f^*(\lambda^*) < f^*(\lambda)$ for any $\lambda \in \Delta^{N-1}, \lambda \neq \lambda^*$.*

---

[25]For any subset $\mathcal{I} \subset [N]$ and $\mathcal{J} \subset [K]$, write $C_{\mathcal{I}\mathcal{J}}$ for the sub-matrix of $C$ with row indices in $\mathcal{I}$ and column indices in $\mathcal{J}$. Likewise, let $C_{-\mathcal{I}\mathcal{J}}$ be the sub-matrix of $C$ after deleting rows in $\mathcal{I}$ and columns in $\mathcal{J}$.

*Proof.* First, we will assume that $[(C^*)^{-1}]_{1i}$ is positive for $1 \leq i \leq K$. This is without loss because we can always work with the "negative" of any signal (replace $c_i$ with $-c_i$), which does not affect agents' behavior.

Since $f(q_1, \ldots, q_N)$ is convex in its arguments, $f^*(\lambda) = \lim_{t \to \infty} t \cdot f(\lambda_1 t, \ldots, \lambda_N t)$ is also convex in $\lambda$. To show $f^*(\lambda^*) < f^*(\lambda)$, we only need to show $f^*(\lambda^*) < f^*((1-\epsilon)\lambda^* + \epsilon\lambda)$ for some $\epsilon > 0$. In other words, it suffices to show $f^*(\lambda^*) < f^*(\lambda)$ for $\lambda$ in an $\epsilon$-neighborhood of $\lambda^*$. By assumption, $\mathcal{S}^*$ is minimally-spanning and so its signals are linearly independent. Thus its signals must span all of the $K$ states. From this it follows that the $K \times K$ matrix $C'\Lambda^*C$ is positive definite, and by (7) the function $f^*$ is differentiable near $\lambda^*$ (not just , see Remark 1 below).

We claim that the partial derivatives of $f^*$ satisfy the following inequality:

$$\partial_K f^*(\lambda^*) < \partial_i f^*(\lambda^*) \leq 0, \forall i > K. \tag{**}$$

Once this is proved, we will have, for $\lambda$ close to $\lambda^*$,

$$f^*(\lambda_1, \ldots, \lambda_K, \lambda_{K+1}, \ldots, \lambda_N) \geq f^*(\lambda_1, \ldots, \lambda_{K-1}, \lambda_K + \lambda_{K+1} + \cdots + \lambda_N, 0, \ldots, 0) \geq f^*(\lambda^*). \tag{8}$$

The first inequality is based on (**) and continuous differentiability of $f^*$, while the second inequality is because $\lambda^*$ uniquely minimizes $f^*$ if society only observes the first $K$ signals. Moreover, when $\lambda \neq \lambda^*$, one of these inequalities is strict so that $f^*(\lambda) > f^*(\lambda^*)$ strictly.

To prove (**), we recall that

$$f^*(\lambda_1, \ldots, \lambda_N) = e_1'(C'\Lambda C)^{-1}e_1.$$

Since $\Lambda = \text{diag}(\lambda_1, \ldots, \lambda_N)$, its derivative is $\partial_i \Lambda = \Delta_{ii}$, which is an $N \times N$ matrix whose $(i,i)$-th entry is 1 and all other entries are zero. Using properties of matrix derivatives, we obtain

$$\partial_i f^*(\lambda) = -e_1'(C'\Lambda C)^{-1}C'\Delta_{ii}C(C'\Lambda C)^{-1}e_1.$$

As the $i$-th row vector of $C$ is $c_i'$, $C'\Delta_{ii}C$ is the $K \times K$ matrix $c_ic_i'$. The above simplifies to

$$\partial_i f^*(\lambda) = -[e_1'(C'\Lambda C)^{-1}c_i]^2.$$

At $\lambda = \lambda^*$, the matrix $C'\Lambda C$ further simplifies to $(C^*)' \cdot \text{diag}(\lambda_1^*, \ldots, \lambda_K^*) \cdot (C^*)$, which is a product of $K \times K$ invertible matrices. We thus deduce that

$$\partial_i f^*(\lambda^*) = -\left[e_1' \cdot (C^*)^{-1} \cdot \text{diag}\left(\frac{1}{\lambda_1^*}, \ldots, \frac{1}{\lambda_K^*}\right) \cdot ((C^*)')^{-1} \cdot c_i\right]^2.$$

It is crucial for our analysis that the term in the brackets is a linear function of $c_i$. To ease notation, we write $v' = e_1' \cdot (C^*)^{-1} \cdot \text{diag}\left(\frac{1}{\lambda_1^*}, \ldots, \frac{1}{\lambda_K^*}\right) \cdot ((C^*)')^{-1}$ and $\gamma_i = \langle v, c_i \rangle$. Then

$$\partial_i f = -\gamma_i^2, \ 1 \leq i \leq N. \tag{9}$$

For $1 \leq i \leq K$, $((C^*)')^{-1} \cdot c_i$ is just $e_i$. Thus, using the assumption $[(C^*)^{-1}]_{1j} > 0, \forall j$, we have

$$\gamma_i = e_1' \cdot (C^*)^{-1} \cdot \text{diag}\left(\frac{1}{\lambda_1^*}, \ldots, \frac{1}{\lambda_K^*}\right) \cdot e_i = \frac{[(C^*)^{-1}]_{1i}}{\lambda_i^*} = \sum_{j=1}^{K}[(C^*)^{-1}]_{1j} = \phi(\mathcal{S}^*), \ 1 \leq i \leq K. \tag{10}$$

33

On the other hand, choosing any $i > K$, we can uniquely write the vector $c_i$ as a linear combination of $c_1, \ldots, c_K$. By Lemma 5, for any $i > K$ we have

$$\gamma_i = \langle v, c_i \rangle = \sum_{j=1}^{K} \alpha_j \cdot \langle v, c_j \rangle = \sum_{j=1}^{K} \alpha_j \cdot \gamma_j = \phi(\mathcal{S}^*) \cdot \sum_{j=1}^{K} \alpha_j. \tag{11}$$

The last equality uses (10). Since $|\sum_{j=1}^{K} \alpha_j| < 1$, the absolute value of $\gamma_i$ for any $i > K$ is strictly smaller than the absolute value of $\gamma_K$. This together with (9) proves the desired inequality (\*\*), and the lemma follows. $\qquad\square$

**Remark 1.** *The essence of this proof is the following non-obvious fact: the subset $\{1, \ldots, K\}$ uniquely minimizes $\phi$ among all subsets of size $K$ if and only if*

$$\phi(\{1, \ldots, K\}) < \phi(\{1, \ldots, K\} \cup \{i\} \backslash \{j\}), \ \ \forall 1 \leq j \leq K < i \leq N.$$

*That is, if a set of $K$ signals does not minimize $\phi$, then we can improve the speed of learning simply by adding* one *signal to replace* one *existing signal. This property enables us to reduce the general problem with $N$ signals to the much simpler problem with $K+1$ signals, and we are able to use calculus to resolve the latter problem, see (\*\*).*

*However, the above fact relies on the original set containing exactly $K$ signals. To see this, consider two states and three signals with coefficient vectors $c_1 = (0.5, 0), c_2 = (1, 1), c_3 = (1, -1)$. If we start with the first signal alone, adding either of the latter two signals does not decrease $\phi$. However, the latter two signals combined yield a faster speed of learning, as $\phi(\{2, 3\}) = 1 < 2 = \phi(\{1\})$. On the technical level, this occurs because at $\lambda = (1, 0, 0)$, $f^*$ is differentiable along every direction but not differentiable as a multivariate function (i.e. it does not admit a gradient vector). Thus, even though the partial derivatives satisfy (\*\*), we cannot deduce that any* directional *derivative similarly satisfies (\*\*). It is for this reason that we need a different proof of Lemma 6 when $|\mathcal{S}^*| < K$, which we present later.*

**Remark 2.** *Still assuming that the "best" subset $\mathcal{S}^*$ contains exactly $K$ signals, we now show $n_i(t) = \lambda_i^* \cdot t + O(1), \forall i$, which improves upon the conclusion of Theorem 1. First, we can apply Lemma 5 to find a positive constant $\eta < 1$ such that for each $i > K$, if $c_i = \sum_{j=1}^{K} \alpha_j c_j$ then $|\sum_{j=1}^{K} \alpha_j| \leq 1 - \eta$. By (9), (10) and (11), we have*

$$\partial_1 f(\lambda^*) = \cdots = \partial_K f(\lambda^*) = -\phi(\mathcal{S}^*)^2; \quad \partial_i f(\lambda^*) \geq -(1 - \eta)^2 \cdot \phi(\mathcal{S}^*)^2, \forall i > K. \tag{12}$$

For any $\lambda \in \Delta^{N-1}$, the convexity of $f^*$ implies[26]

$$f^*(\lambda) \geq f^*(\lambda^*) + \sum_{i=1}^{N}(\lambda_i - \lambda_i^*) \cdot \partial_i f^*(\lambda^*)$$

$$= f^*(\lambda^*) + \sum_{i=1}^{N}(\lambda_i - \lambda_i^*) \cdot (\partial_i f^*(\lambda^*) + \phi(\mathcal{S}^*)^2) \tag{13}$$

$$\geq f^*(\lambda^*) + (2\eta - \eta^2) \cdot \phi(\mathcal{S}^*)^2 \cdot \sum_{i=K+1}^{N} \lambda_i.$$

The second line uses $\sum_{i=1}^{N}(\lambda_i - \lambda_i^*) = 0$ and the last inequality is due to *(12)*.

Consider any division $(q_1, \ldots, q_N)$ at time $t$. A straightforward refinement of Lemma *3* gives that whenever $f^*(\lambda)$ is finite, $t \cdot f(\lambda t)$ approaches $f^*(\lambda)$ at the rate of $\frac{1}{t}$. In particular $f(\lambda^* \cdot t) = \frac{1}{t} \cdot f^*(\lambda^*) + O(\frac{1}{t^2})$. For $(q_1, \ldots, q_N)$ to be a $t$-optimal division, it is necessary that $f(q_1, \ldots, q_N) \leq f(\lambda^* \cdot t)$. Thus

$$f^*\left(\frac{q_1}{t}, \ldots, \frac{q_N}{t}\right) \leq f^*(\lambda^*) + O\left(\frac{1}{t}\right). \tag{14}$$

By *(13)* and *(14)*, any $t$-optimal division $n(t)$ must satisfy $n_i(t) = O(1)$ for each signal $i > K$. Conditional on these signal counts, society's optimal choice over signals 1 through $K$ must satisfy $n_i(t) = \lambda_i^* \cdot t + O(1), \forall 1 \leq i \leq K$, as shown in Proposition *1*. This is what we desire to prove here.

## C.5 A Perturbation Argument

We have shown that whenever $\phi(\mathcal{S})$ is uniquely minimized by a set $\mathcal{S}$ containing $K$ signals,

$$\min_{\lambda \in \Delta^{N-1}} f^*(\lambda) = f^*(\lambda^*) = \min_{\mathcal{S} \subset [N]} \phi(\mathcal{S})^2 = \phi([N])^2$$

We now show this equality holds more generally.

**Lemma 7.** *For any coefficient matrix $C$,*

$$\min_{\lambda \in \Delta^{N-1}} f^*(\lambda) = \phi([N])^2. \tag{15}$$

*Proof.* Because society can choose to focus on any minimal spanning set, it is clear that $\min_\lambda f^*(\lambda) \leq \phi([N])^2 = \min_{\mathcal{S}}(\phi(\mathcal{S}))^2$. It remains to prove $f^*(\lambda) \geq \phi([N])^2$ for any fixed $\lambda \in \Delta^{N-1}$. By Lemma *3*, we need to show $[(C'\Lambda C)^{-1}]_{11} \geq \phi([N])^2$.

---

[26] As mentioned in Remark 1, it is crucial that $f^*$ is differentiable at $\lambda^*$. The argument here relies on the directional derivative in the direction $\lambda - \lambda^*$ being well-defined and equal to a linear sum of partial derivatives.

This was already proved for *generic* coefficient matrices $C$; specifically, those for which $\phi(\mathcal{S})$ is minimized by a set of $K$ signals. But even if $C$ is "non-generic", we can approximate it by a sequence of "generic" matrices $C_m$.[27] Along this sequence, we have

$$[(C_m'\Lambda C_m)^{-1}]_{11} \geq \phi_m([N])^2$$

where $\phi_m$ is the speed of learning from the $N$ signals given by coefficient matrix $C_m$. As $m \to \infty$, the LHS above approaches $[(C'\Lambda C)^{-1}]_{11}$. Thus the lemma will follow once we show that $\limsup_{m\to\infty} \phi_m([N]) \geq \phi([N])$.

For this we invoke the following characterization

$$\phi([N]) = \min_{\beta \in \mathbb{R}^N} \sum_{i=1}^{N} |\beta_i| \quad \text{s.t.} \quad e_1 = \sum_{i=1}^{N} \beta_i \cdot c_i.$$

If $e_1 = \sum_i \beta_i^{(m)} \cdot c_i^{(m)}$ along the convergent sequence, then $e_1 = \sum_i \beta_i \cdot c_i$ for any limit point $\beta$ of $\beta^{(m)}$. This enables us to conclude $\liminf_{m\to\infty} \phi_m([N]) \geq \phi([N])$, which is more than what we need. $\qquad\square$

## C.6 Proof of Theorem 1 when $|\mathcal{S}^*| < K$

We now complete the proof of Theorem 1 for the case where the "best" subset $\mathcal{S}^*$ contains fewer than $K$ signals. To be precise, let $\mathcal{S}^* = \{1, \ldots, k\}$ and define $\lambda^* \in \Delta^{N-1}$ to be the optimal frequencies when only the first $k$ signals are observed. We will show $n_i(t) \sim \lambda_i^* \cdot t, \forall i$. By Lemma 4, we only need to show that $\lambda^*$ uniquely minimizes $f^*(\lambda)$ over the simplex. Since $f^*(\lambda^*) = \phi(\mathcal{S}^*)^2 = \phi([N])^2$ by definition, we know from Lemma 7 that $\lambda^*$ does minimize $f^*(\lambda)$.

It remains to show that $\lambda^*$ is the *unique* minimizer. Suppose for contradiction that $f^*(\lambda^*) = f^*(\tilde{\lambda})$ for some $\tilde{\lambda} \in \Delta^{N-1}$ distinct from $\lambda^*$. For $\eta \in \mathbb{R}$, define $\lambda^\eta = \lambda^* + \eta \cdot (\tilde{\lambda} - \lambda^*)$, so that $\lambda^0 = \lambda^*, \lambda^1 = \tilde{\lambda}$. Observe that when $\eta \in (0, 1)$, $\lambda^\eta$ is a convex combination between $\lambda^*$ and $\tilde{\lambda}$. Thus the convexity of $f^*$ implies

$$f^*(\lambda^\eta) \leq (1 - \eta)f^*(\lambda^*) + \eta f^*(\tilde{\lambda}) = f^*(\lambda^*)$$

Since $f^*(\lambda^*)$ is minimal, we must then have $f^*(\lambda^\eta) = f^*(\lambda^*)$ for $\eta \in (0, 1)$. But for fixed $\lambda^*$ and $\lambda$, (7) shows that the value of $f^*(\lambda^\eta)$ is a rational function (quotient of two polynomials) of $\eta$. Thus this rational function is itself a constant. Consequently, $f^*(\lambda^\eta) = f^*(\lambda^*)$ for all $\eta$ (not just those in the unit interval) such that $\lambda^\eta \in \Delta^{N-1}$.

Because $\tilde{\lambda} \neq \lambda^*$, there exists some $j \in \{1, \ldots, k\}$ such that $\tilde{\lambda}_j < \lambda_j^*$. Without loss, we assume $\frac{\tilde{\lambda}_1}{\lambda_1^*}$ is the smallest among such ratios. Let $\eta = \frac{\lambda_1^*}{\lambda_1^* - \tilde{\lambda}_1}$, then the vector $\lambda^\eta$

---

[27]First, we may add repetitive signals to ensure $N \geq K$. This does not affect the value of $\min f^*(\lambda)$ or $\phi([N])$. Whenever $N \geq K$, it is generically true that every minimal spanning set contains exactly $K$ signals. Moreover, the equality $\phi(\mathcal{S}) = \phi(\tilde{\mathcal{S}})$ for $\mathcal{S} \neq \tilde{\mathcal{S}}$ induces a non-trivial polynomial equation over the entries in $C$. This means we can always find $C^{(m)}$ close to $C$ such that for the coefficient matrix $C^{(m)}$, different subsets $\mathcal{S}$ (of size $K$) attain different values of $\phi(\mathcal{S})$.

has first-coordinate 0 and all other coordinates non-negative. By our preceding analysis, $f^*(\lambda^\eta) = f^*(\lambda^*)$ for this $\eta$. However, since $\lambda^\eta$ "ignores" signal 1, Lemma 7 implies that

$$f^*(\lambda^\eta) \geq \min_{\lambda \in \Delta^{N-1}, \ \lambda_1 = 0} f^*(\lambda) = \phi([N]\backslash\{1\})^2.$$

By assumption, $\mathcal{S}^* = \{1, \ldots, k\}$ is the *unique* minimal spanning set that minimizes $\phi$. Thus the RHS above is strictly larger than $\phi(\mathcal{S}^*)^2 = f^*(\lambda^*)$, leading to the contradictory result $f^*(\lambda^\eta) > f^*(\lambda^*)$.

This contradiction shows $\lambda^*$ must uniquely minimize $f^*(\lambda)$. Theorem 1 follows.

# D    Proof of Theorem 2

Let signals $1, \ldots, k$ (with $k \leq K$) be a minimally spanning set that is optimal in its subspace. We will demonstrate an open set of prior beliefs given which *all agents* observe these $k$ signals. Since these signals are minimally spanning, they must be linearly independent. Thus we can consider linearly transformed states $\tilde{\theta}_1, \ldots, \tilde{\theta}_K$ such that these $k$ signals are simply $\tilde{\theta}_1, \ldots, \tilde{\theta}_k$ plus standard Gaussian noise. This linear transformation is invertible, so any prior over the original states is bijectively mapped to a prior over the transformed states. Thus it is without loss to work with the transformed model and look for prior beliefs over the transformed states.

The payoff-relevant state $\omega$ becomes a linear combination $w_1\tilde{\theta}_1 + \cdots + w_k\tilde{\theta}_k$. We may without loss assume the weights $w_i$ are all positive. Moreover, since the first $k$ signals are optimal in its subspace, Lemma 5 implies that any other signal that belongs to this subspace can be written as

$$\sum_{i=1}^{k} \alpha_i \tilde{\theta}_i \ + \ \mathcal{N}(0,1)$$

with $|\sum_{i=1}^{k} \alpha_i| < 1$. On the other hand, if a signal does not belong to this subspace, it must take the form of

$$\sum_{i=1}^{K} \beta_i \tilde{\theta}_i \ + \ \mathcal{N}(0,1)$$

with $\beta_{k+1}, \ldots, \beta_K$ not all equal to zero.

Now consider a prior belief such that $\tilde{\theta}_1, \ldots, \tilde{\theta}_K$ are independent from each other. Given prior variances $v_1, \ldots, v_K$, the reduction in the variance of $w_1\tilde{\theta}_1 + \cdots + w_k\tilde{\theta}_k$ by any signal $\sum_{i=1}^{k} \alpha_i\tilde{\theta}_i + \mathcal{N}(0,1)$ is

$$\frac{(\sum_{i=1}^{k} \alpha_i w_i v_i)^2}{1 + \sum_{i=1}^{k} \alpha_i^2 v_i}$$

If $v_1, \ldots, v_k$ are small positive numbers and if the product $w_i v_i$ is approximately constant across $1 \leq i \leq k$, then the above is approximately $(\sum_{i=1}^{k} \alpha_i)^2 w_1^2 v_1^2$. Since $|\sum_{i=1}^{k} \alpha_i| < 1$, we deduce that any other signal belonging to the subspace of the first $k$ signals is worse than signal 1 (in the first period), whose variance reduction is $\frac{w_1^2 v_1^2}{1+v_1}$.

Meanwhile, take any signal that does not belong to the subspace. The variance reduction by such a signal $\sum_{i=1}^{K} \beta_i \tilde{\theta}_i + \mathcal{N}(0,1)$ is

$$\frac{(\sum_{i=1}^{k} \beta_i w_i v_i)^2}{1 + \sum_{i=1}^{K} \beta_i^2 v_i}$$

As $\beta_{k+1}, \ldots, \beta_K$ are not all zero, the denominator above can be arbitrarily large if $v_{k+1}, \ldots, v_K$ are chosen to be large. Then, this signal is again worse than signal 1 for the first agent, similar to the situation in Example 4.

To summarize, we have shown that whenever the prior variances $v_1, \ldots, v_K$ satisfy the following three conditions, the first agent chooses among the first $k$ signals:

1. $v_1, \ldots, v_k$ are close to 0;

2. $w_1 v_1, \ldots, w_k v_k$ have pairwise ratios close to 1;

3. $v_{k+1}, \ldots, v_K$ are large.[28]

To show that *every agent* chooses among the first $k$ signals, it suffices to check that starting from any prior beliefs satisfying the above conditions, the posterior beliefs after observing a signal continue to satisfy these conditions. Since variances decrease over time, the first condition is obviously satisfied. By independence, learning about $\tilde{\theta}_1, \ldots, \tilde{\theta}_k$ does not affect the variances of the remaining states. So $v_{k+1}, \ldots, v_K$ are unchanged, and the third condition is verified. Finally, the second condition holds for the posterior beliefs because the signal $i$ that is chosen has the greatest value of $\frac{w_i^2 v_i^2}{1 + v_i}$. This choice ensures that $v_i \propto \frac{1}{w_i}$, as shown also in Liang, Mu and Syrgkanis (2017). Hence Theorem 2 is proved.

Strictly speaking, the above construction does not provide an *open set* of prior beliefs given which agents always observe the first $k$ signals. This is because we restricted attention to priors that are independent over $\tilde{\theta}_1, \ldots, \tilde{\theta}_K$. But it could be shown that the argument extends to mild correlation across states. We omit the somewhat cumbersome details, which do not add any further intuition.

# E    Proof of Theorem 3

## E.1    Preliminary Steps

Given any prior, let $\mathcal{A} \subset [N]$ be the set of signals that are observed by infinitely many agents. We first show that $\mathcal{A}$ is a spanning set.

Indeed, by definition we can find some period $t$ after which agents only observe signals in $\mathcal{A}$. Also note that the variance reduction of any signal approaches zero as its signal count gets large. Thus, along society's signal path, the variance reduction is close to zero at sufficiently late periods.

---

[28]Formally, we require that for some $\xi > 0$, it holds that $v_1, \ldots, v_k < \xi$; $\max_{1 \le i \le k} w_i v_i \le (1 + \xi) \cdot \min_{1 \le i \le k} w_i v_i$; and $v_{k+1}, \ldots, v_K > \frac{1}{\xi}$.

If $\mathcal{A}$ is not spanning, society's posterior variance remains bounded away from zero. Thus in the limit where each signal in $\mathcal{A}$ has infinite signal counts, there still exists some signal $j$ outside of $\mathcal{A}$ whose variance reduction is strictly positive.[29] By continuity, at sufficiently late periods, observing signal $j$ would reduce the variance by a positive amount. This is a profitable deviation from observing some signal in $\mathcal{A}$, leading to a contradiction!

Now that $\mathcal{A}$ is spanning, we can take $\mathcal{S}$ to be the optimal minimal spanning set in the subspace spanned by $\mathcal{A}$. To prove Theorem 3, we will show the long-run frequencies are positive precisely for the signals in $\mathcal{S}$. Ignoring the initial periods, it is without loss to assume that only signals in $\overline{\mathcal{A}}$ are available. It suffices to show that whenever the signals observed infinitely often *span that subspace*, agents eventually sample from the optimal subset $\mathcal{S}$. To ease notation, we assume this subspace is the entire $\mathbb{R}^K$, and prove the following result:

**Theorem 3 Restated.** Suppose that the signals observed infinitely often span $\mathbb{R}^K$. Then society eventually observes signals in $\mathcal{S}^*$ with frequencies $\lambda^*$.

The next sections are devoted to the proof of this restatement.

## E.2  Controlling the Derivatives

To study the posterior variance function $f$, it will be convenient to instead work with the homogenous function $f^*$ we introduced in Lemma 3. We formalize this connection as follows:

**Lemma 8.** *Suppose that signals in $\mathcal{A}$ span $\mathbb{R}^K$. Then, as $q_i \to \infty$ for each $i \in \mathcal{A}$,*

$$f(q_1, \ldots, q_N) \sim \frac{1}{t} \cdot f^* \left( \frac{q_1}{t}, \ldots, \frac{q_N}{t} \right) \quad with \quad t = \sum_{i=1}^N q_i$$

*The partial derivatives and second partial derivatives also satisfy the approximations*

$$\partial_j f(q_1, \ldots, q_N) \sim \frac{1}{t^2} \cdot \partial_j f^* \left( \frac{q_1}{t}, \ldots, \frac{q_N}{t} \right)$$

$$\partial_{jj} f(q_1, \ldots, q_N) \sim \frac{1}{t^3} \cdot \partial_{jj} f^* \left( \frac{q_1}{t}, \ldots, \frac{q_N}{t} \right)$$

*Proof.* Recall that

$$f(q_1, \ldots, q_N) = \left[ ((V^0)^{-1} + C'QC)^{-1} \right]_{11}.$$

Since $q_i \to \infty$ for $i \in \mathcal{A}$, the least eigenvalue of the matrix $C'QC$ approaches infinity. That is, for any $\epsilon > 0$, it holds eventually that $(V^0)^{-1} \preceq \epsilon \cdot C'QC$ in matrix order. Then

$$\frac{1}{1+\epsilon} \cdot [(C'QC)^{-1}]_{11} \leq f(q_1, \ldots, q_N) \leq [(C'QC)^{-1}]_{11}.$$

---

[29]Formally, let $s_1, \ldots, s_N$ denote the limit signal counts, where $s_i = \infty$ if and only if $i \in \mathcal{A}$. Then there exists $j$ such that $f(s_j + 1, s_{-j}) < f(s_j, s_{-j})$. This is because if $f(s_j + 1, s_{-j}) = f(s_j, s_{-j})$ for each $j$, then the partial derivatives of $f$ at $s$ are all zero. Since $f$ is differentiable, this would imply all directional derivatives of $f$ are also zero. By the convexity of $f$, $f(s)$ must achieve minimum value. But by assumption there exists a spanning set, so $f(q) = 0$ if $q_1, \ldots, q_N$ are all infinite. This contradicts $f(s) > 0$.

Equivalently, this shows

$$\frac{1}{(1+\epsilon)t} \cdot f^* \left( \frac{q_1}{t}, \ldots, \frac{q_N}{t} \right) \leq f(q_1, \ldots, q_N) \leq \frac{1}{t} \cdot f^* \left( \frac{q_1}{t}, \ldots, \frac{q_N}{t} \right).$$

Similar approximation holds for the derivatives, proving the lemma. $\qquad\square$

**Lemma 9.** *Under the same assumptions as in Lemma 8, it holds that*

$$\frac{\partial_{jj} f(q_1, \ldots, q_N)}{\partial_j f(q_1, \ldots, q_N)} \to 0$$

*and similarly*

$$\frac{\partial_{jj} f^* \left( \frac{q_1}{t}, \ldots, \frac{q_N}{t} \right)}{t \cdot \partial_j f^* \left( \frac{q_1}{t}, \ldots, \frac{q_N}{t} \right)} \to 0$$

*Proof.* It suffices to prove the first result. From $f(q_1, \ldots, q_N) = e_1' \cdot [(V^0)^{-1} + C'QC]^{-1} \cdot e_1$ we compute that

$$\partial_j f = -e_1' \cdot [(V^0)^{-1} + C'QC]^{-1} \cdot c_j \cdot c_j' \cdot [(V^0)^{-1} + C'QC]^{-1} \cdot e_1$$

and

$$\partial_{jj} f = 2e_1' \cdot [(V^0)^{-1} + C'QC]^{-1} \cdot c_j \cdot c_j' \cdot [(V^0)^{-1} + C'QC]^{-1} \cdot c_j \cdot c_j' \cdot [(V^0)^{-1} + C'QC]^{-1} \cdot e_1.$$

Let $\gamma_j = e_1' \cdot [(V^0)^{-1} + C'QC]^{-1} \cdot c_j$, which is a number. Then the above shows

$$\partial_j f = -\gamma_j^2; \qquad \partial_{jj} f = 2\gamma_j^2 \cdot c_j' \cdot [(V^0)^{-1} + C'QC]^{-1} \cdot c_j.$$

Again, all eigenvalues of the matrix $(V^0)^{-1} + C'QC$ become large as $q_i \to \infty$ for $i \in \mathcal{A}$. Thus for arbitrarily large constant $L$, eventually $(V^0)^{-1} + C'QC \succeq L \cdot c_j c_j'$ in matrix norm. Then the number $c_j' \cdot [(V^0)^{-1} + C'QC]^{-1} \cdot c_j$ is arbitrarily small, and the above display shows $\partial_{jj} f$ is small compared to $\partial_j f$. $\qquad\square$

The above lemmata imply that at sufficiently late periods along society's signal path, the variance reduction of any *discrete* signal can be approximated by the continuous partial derivative of $f$ (or $f^*$). A direct corollary is the following:

**Lemma 10.** *For any $\epsilon > 0$, there exists sufficiently large $t(\epsilon)$ such that if signal $j$ is observed in any period $t + 1$ later than $t(\epsilon)$, then*

$$\partial_j f^* \left( \frac{m(t)}{t} \right) \leq (1 - \epsilon) \min_{1 \leq l \leq N} \partial_l f^* \left( \frac{m(t)}{t} \right).$$

That is, the signal choice in any sufficiently late period *almost* minimizes the directional derivative of $f^*$.

## E.3   (Pseudo) Gradient Descent of $f^*$

We define $\lambda(t) = \frac{m(t)}{t} \in \Delta^{N-1}$. If $j$ is the signal choice in period $t+1$, then it is easily checked that

$$\lambda(t+1) = \frac{t}{t+1}\lambda(t) + \frac{1}{t+1}e_j.$$

The frequencies $\lambda(t)$ move in the direction of $e_j$, which is the direction where $f^*$ decreases almost the fastest (by Lemma 10). Thus, the evolution of $\lambda(t)$ over time resembles the gradient descent dynamics—the value of $f^*(\lambda(t))$ roughly decreases over time, and we can expect that eventually $\lambda(t)$ approaches the unique minimizer $\lambda^*$ of $f^*$.

To formalize this intuition, we consider (for fixed $\epsilon > 0$ and sufficiently large $t$)

$$
\begin{aligned}
f^*(\lambda(t+1)) &= f^*\left(\frac{t}{t+1}\lambda(t) + \frac{1}{t+1}e_j\right)\\
&= f^*\left(\frac{t}{t+1}\lambda(t)\right) + \frac{1}{t+1}\cdot\partial_j f^*\left(\frac{t}{t+1}\lambda(t)\right) + O\left(\frac{1}{(t+1)^2}\cdot\partial_{jj} f^*\left(\frac{t}{t+1}\lambda(t)\right)\right)\\
&\leq f^*\left(\frac{t}{t+1}\lambda(t)\right) + \frac{1-\epsilon}{t+1}\cdot\partial_j f^*\left(\frac{t}{t+1}\lambda(t)\right)\\
&= \frac{t+1}{t}\cdot f^*(\lambda(t)) + \frac{(1-\epsilon)(t+1)}{t^2}\cdot\partial_j f^*(\lambda(t))\\
&\leq f^*(\lambda(t)) + \frac{1}{t}\cdot f^*(\lambda(t)) + \frac{1-2\epsilon}{t}\cdot\min_{1\leq l\leq N}\partial_l f^*(\lambda(t)).
\end{aligned}
\tag{16}
$$

The first inequality uses Lemma 9, the next equality uses the homogeneity of $f^*$, and the last inequality uses Lemma 10.

Write $\lambda = \lambda(t)$ for short. Observe that $f^*$ is differentiable at $\lambda$, since $\lambda_i(t) > 0$ for $i \in \mathcal{A}$, which spans the entire space. Thus the convexity of $f^*$ yields

$$f^*(\lambda^*) \geq f^*(\lambda) + \sum_{j=1}^{N}(\lambda_j^* - \lambda_j)\cdot\partial_j f^*(\lambda).$$

The homogeneity of $f^*$ implies $\sum_{j=1}^{N}\lambda_j\cdot\partial_j f^*(\lambda) = -f^*(\lambda)$. This enables us to rewrite the above display as

$$\sum_{j=1}^{N}\lambda_j^*\cdot\partial_j f^*(\lambda) \leq f^*(\lambda^*) - 2f^*(\lambda).$$

Thus, in particular,

$$\min_{1\leq l\leq N}\partial_l f^*(\lambda(t)) \leq f^*(\lambda^*) - 2f^*(\lambda). \tag{17}$$

Combining (16) and (17), we have for all large $t$:

$$f^*(\lambda(t+1)) \leq f^*(\lambda(t)) + \frac{1}{t}\cdot[(1-2\epsilon)\cdot f^*(\lambda^*) - (1-4\epsilon)\cdot f^*(\lambda(t))]. \tag{18}$$

We claim this implies $f^*(\lambda(t)) \leq (1+4\epsilon)\cdot f^*(\lambda^*)$ holds for all large $t$. Indeed, if this holds for *some* $t$, then (18) implies the same is true at future periods. It thus suffices to show

the opposite inequality $f^*(\lambda(t)) > (1 + 4\epsilon) \cdot f^*(\lambda^*)$ cannot hold at every large $t$. By (18), that would give $f^*(\lambda(t+1)) \leq f^*(\lambda(t)) - \frac{\epsilon \cdot f^*(\lambda^*)}{t}$. But since the harmonic series diverges, $f^*(\lambda(t))$ would then decrease without bound, leading to a contradiction!

Hence we have shown that for any fixed $\epsilon$, $f^*(\lambda(t)) \leq (1 + 4\epsilon) \cdot f^*(\lambda^*)$ holds eventually. As $\lambda^*$ is the unique minimizer of $f^*$, this implies $\lambda(t) \to \lambda^*$. Theorem 3 follows.

**Remark 3.** *The above argument leaves open the possibility that some signals outside of $\mathcal{S}^*$ are observed* infinitely often, *yet with* zero long-run frequency. *We conjecture this cannot happen, but we are only able to show this when $|\mathcal{S}^*| = K$.*

*Specifically, suppose $|\mathcal{S}^*| = K$ and $m_i(t) \sim \lambda_i^* \cdot t, \forall i$, then we claim that the stronger conclusion $m_i(t) = \lambda_i^* \cdot t + O(1)$ also holds.*[30] *Together with Remark 2, this suggests that the difference between $m_i(t)$ and the optimal $n_i(t)$ remains bounded.*

*To prove this claim, we assume without loss that $\mathcal{S}^* = \{1, \ldots, K\}$ is the first $K$ signals. By the previously established (\*\*), the first $K$ partial derivatives of $f^*$ are equal at $\lambda^*$ and they are strictly smaller (i.e., more negative) than the other partial derivatives. Since these partial derivatives are continuous, we can find $\epsilon > 0$ such that whenever $\lambda$ is within $\epsilon$ distance from $\lambda^*$, it holds that*

$$\partial_i f^*(\lambda) < (1 + \epsilon) \cdot \partial_j f^*(\lambda), \quad \forall 1 \leq i \leq K < j$$

*By assumption we have $\lambda(t) = \frac{m(t)}{t} \to \lambda^*$. Thus at sufficiently late periods, Lemma 10 implies that the signal choice must be within the first $K$ signals. This shows signals outside of $\mathcal{S}^*$ are observed finitely often, as desired. And for any signal $i$ in $\mathcal{S}^*$, its signal count satisfies $m_i(t) = \lambda_i^* \cdot t + O(1)$ by Proposition 1.*

# F    Proof of Proposition 2

We will prove that given any prior belief, the planner can provide $K - 1$ sufficiently precise signals so that once they are processed, society eventually observes the best set $\mathcal{S}^*$. In fact, the following argument shows that the planner can provide these free signals at any time $t$, not necessarily before agents arrive.

The proof of the proposition closely resembles the proof of the restated Theorem 3, see Appendix E. Indeed, with sufficiently high precision on the free signals, it is as if each free signal has unit precision but is observed many times. Thus, as long as the $K - 1$ free signals span $b_1, \ldots, b_{K-1}$, the restated Theorem 3 applies since society eventually learns $\omega$ anyways. Of course, the assumption of that theorem is not exactly satisfied, and one may wonder whether *observing a signal many times has the same consequence as observing it infinitely often.* In what follows we show how to resolve this concern.

Consider for simplicity that the planner provides $L$ i.i.d. free signals in $\mathcal{A} \subset [N]$, which spans $b_1, \ldots, b_{K-1}$. We are free to choose $L$ by making $\gamma$ sufficiently large. Then, at any time $t$, the signal count $m_i(t)$ is at least $L$ for each signal $i \in \mathcal{A}$. Fix any $\epsilon > 0$, there exists such an $L$ that the approximations in Lemma 8 and 9 hold up to a margin of error no more

---

[30]Thus, the conclusion of Corollary 5 can be strengthened.

than $\epsilon$. That is, for Lemma 8, we now have

$$(1 - \epsilon) \cdot f(q_1, \ldots, q_N) \leq \frac{1}{t} \cdot f^* \left( \frac{q_1}{t}, \ldots, \frac{q_N}{t} \right) \leq (1 + \epsilon) \cdot f(q_1, \ldots, q_N)$$

etc., and we similarly modify Lemma 9 to

$$|\frac{\partial_{jj} f(q_1, \ldots, q_N)}{\partial_j f(q_1, \ldots, q_N)}| \leq \epsilon.$$

These hold because the signal precision matrix $C'QC$ eventually dominates the prior precision matrix $(V^0)^{-1}$.

As a result, for any fixed $\epsilon$, Lemma 10 still holds if we choose $L$ to be sufficiently large. We could then derive (16), (17) and (18) in the same way as before. This enables us to conclude

$$f^*(\lambda(t)) \leq (1 + 4\epsilon) \cdot f^*(\lambda^*)$$

at every late period $t$. Since we have fixed $\epsilon$ (and $L$), the above inequality does not by itself imply $\lambda(t) \to \lambda^*$. However, if we had chosen $\epsilon$ to be sufficiently small, then $\lambda(t)$ eventually belongs to a small neighborhood of $\lambda^*$. In particular, for $\epsilon$ small the above display implies $\lambda_i(t) \geq \frac{\lambda_i^*}{2} > 0$ for each $i \in \mathcal{S}^*$.

With such a choice of $\epsilon$ and corresponding $L$, we know that society observes each signal in $\mathcal{S}^*$ with positive frequencies. But Theorem 3 shows that the set of signals with positive frequencies is a minimal spanning set. So this set must be $\mathcal{S}^*$ itself, and the long-run frequencies must be $\lambda^*$. Hence efficient learning is achieved, and Proposition 2 follows.

# G   Multiple Payoff-Relevant States

In this appendix, we consider optimal long-run acquisitions for the problem of predicting multiple states. We assume that society seeks to minimize the sum of his posterior variances about the $K$ states. Formally, the planner's objective function is to minimize

$$F(q_1, \ldots, q_N) = Tr \left[ ((V^0)^{-1} + C'QC)^{-1} \right].$$

subject to the signal counts $q_i$ being integers and summing up to $t$. We use "$Tr$" to denote the trace of a matrix.

The solution to this minimization problem turns out to be very complex when $N > K$. To make the problem more tractable, we impose a further assumption that the signal coefficient vectors $c_i$ have the same norm. This allows us to focus the analysis on the directions of the signals, rather than their precisions.

**Assumption 5** (Unit Norm). *Each vector $c_i \in \mathbb{R}^K$ has norm 1.*

Given this assumption, a basic question is to understand how fast society can jointly learn about different states. If $N = K$ and each signal is about an individual state, then obviously society cannot do better (in the long run) than spending the same number ($\frac{t}{K}$) of observations on each signal. In so doing, its posterior variance at time $t$ about each state is approximately $\frac{K}{t}$, and the sum of these variances is $\frac{K^2}{t}$. Our next result shows this is asymptotically best, even when additional signals are available.

**Proposition 3.** *Under Assumption 5, we have*

$$\liminf_{q_1+\cdots+q_N\to\infty}(q_1+\cdots+q_N)\cdot F(q_1,\ldots,q_N)\geq K^2.$$

For the special case of $K = 2$, we are able to determine the exact asymptotic variance (the value of the LHS above) for any given set of signals, see later in this appendix. Deriving the analogous result for general $K$ is left for future work.

We highlight that unlike the case of a single payoff-relevant state, here the minimum asymptotic variance can in general be achieved by more than one vector of frequencies. Thus, the above results only describe agents' payoffs at large $t$, but they do not pin down agents' optimal behavior. When $q_i$ is not restricted to integer values, Chaloner (1984) showed that the minimum posterior variance at any *fixed* time $t$ is achieved by focusing on at most $\frac{K(K+1)}{2}$ signals. However, it is not known whether the same subset of $\frac{K(K+1)}{2}$ signals are observed for all large $t$, and her result also does not apply to our integer design problem.

## G.1 Proof of Proposition 3

We first show that

$$F^*(\lambda) := \lim_{t\to\infty} t\cdot F(\lambda t) = Tr\left[(C'\Lambda C)^{-1}\right] \tag{19}$$

If at least $K$ of $\lambda_1,\ldots,\lambda_N$ are positive, this follows from the previous formula for $F$. Suppose instead that only $\lambda_1,\ldots,\lambda_k$ are positive, with $k < K$. Consider the limit of $Tr\left[(C'\Lambda C)^{-1}\right]$ as $\lambda_{k+1},\ldots,\lambda_N$ approaches zero. In this limit, the $K \times K$ matrix $C'\Lambda C$ approaches a rank $k$ matrix, so an eigenvalue of $C'\Lambda C$ approaches zero. This means an eigenvalue of $(C'\Lambda C)^{-1}$ approaches infinity, and since all its eigenvalues are non-negative by positive-definiteness, we deduce $Tr\left[(C'\Lambda C)^{-1}\right] \to \infty$. Meanwhile, $F(\lambda t)$ is bounded away from zero since the first $k$ signals cannot identify all of the states $\omega,\ldots,\theta_K$. Thus (19) always hold.

We need to show that if each signal coefficient vector $c_i$ has norm 1, then $F^*(\lambda) \geq K^2$ for all $\lambda \in \Delta^{N-1}$. For this, consider the positive-definite $K \times K$ matrix $C'\Lambda C$. Let its $K$ (positive) eigenvalues be $\beta_1,\ldots,\beta_K$, then we have

$$\beta_1 + \cdots + \beta_K = Tr(C'\Lambda C) = \sum_{i=1}^N \lambda_i \sum_{j=1}^K c_{ij}^2 = \sum_{i=1}^N \lambda_i = 1,$$

Observe that the eigenvalues of the inverse matrix $(C'\Lambda C)^{-1}$ are simply $\frac{1}{\beta_1},\ldots,\frac{1}{\beta_K}$. Thus, by (19) and Cauchy-Schwartz inequality,

$$F^*(\lambda) = Tr\left[(C'\Lambda C)^{-1}\right] = \frac{1}{\beta_1} + \cdots + \frac{1}{\beta_K} \geq \frac{K^2}{\beta_1 + \cdots + \beta_K} = K^2.$$

This proves Proposition 3.

## G.2 Characterization of Asymptotic Variance when $K = 2$

Suppose there are just two states and each signal has unit norm, we determine here the exact value of $\min_{\lambda \in \Delta^{N-1}} F^*(\lambda)$ for any given coefficient matrix $C$. By what we have shown, this value (divided by $t$) approximates the minimum of the objective function $F$ that can be achieved given $t$ observations.

Applying Lemma 3 and adding up the variances about $\omega$ and $\theta_2$, we have for $K = 2$,

$$F^*(\lambda) = \frac{\sum_{i=1}^{N} \lambda_i(x_i^2 + y_i^2)}{\sum_{1 \leq i < j \leq N} \lambda_i \lambda_j (x_i y_j - x_j y_i)^2},$$

where each signal coefficient vector $c_i = (x_i, y_i)'$. By Assumption 5, $x_i^2 + y_i^2 = 1$ for each $i$. Thus the numerator above is exactly 1, and we only need to *maximize* the denominator. It will be convenient to parametrize $(x_i, y_i) = (\cos \phi_i, \sin \phi_i)$, with $\phi_i \in [0, \pi)$ distinct from one another.[31] Then, the denominator becomes

$$\sum_{1 \leq i < j \leq N} \lambda_i \lambda_j (x_i y_j - x_j y_i)^2 = \sum_{1 \leq i < j \leq N} \lambda_i \lambda_j \sin^2(\phi_i - \phi_j) = \frac{1}{4} \cdot \sum_{1 \leq i < j \leq N} \lambda_i \lambda_j (2 - 2\cos(2\phi_i - 2\phi_j))$$

$$= \frac{1}{4}(\lambda_1 + \cdots + \lambda_N)^2 - \frac{1}{4}\sum_{i=1}^{N} \lambda_i^2 - \frac{1}{4} \sum_{1 \leq i < j \leq N} \lambda_i \lambda_j 2 \cos(2\phi_i - 2\phi_j)$$

$$= \frac{1}{4} - \sum_{i,j=1}^{N} \lambda_i \lambda_j \cos(2\phi_i - 2\phi_j)$$

$$= \frac{1}{4} - \left(\sum_{i=1}^{N} \lambda_i \cos 2\phi_i\right)^2 - \left(\sum_{i=1}^{N} \lambda_i \sin 2\phi_i\right)^2.$$

This recovers the result of Proposition 3 that $F^*(\lambda) \geq 4$. More generally, given $\phi_1, \ldots, \phi_N$, let $u_i = (\cos 2\phi_i, \sin 2\phi_i)$ be a vector/point lying on the unit circle. Then society seeks to *minimize* $(\sum_{i=1}^{N} \lambda_i \cos 2\phi_i)^2 + (\sum_{i=1}^{N} \lambda_i \sin 2\phi_i)^2$, which is the squared norm of the vector $\sum_{i=1}^{N} \lambda_i u_i$. Taking a geometric perspective, this problem is to choose a point in the convex hull of points $u_1, \ldots, u_N$ that is closest to the origin. There are two possibilities:

1. Suppose the points $u_1, \ldots, u_N$ lie on a semi-circle. Without loss, we label $u_1$ as the point closest to one end of this semi-circle and $u_2$ being closest to the other end. Then the point in $Conv(u_1, \ldots, u_N)$ that is closest to the origin is the mid-point between $u_1$ and $u_2$. In this case $F^*$ is *uniquely* minimized at $\lambda = (\frac{1}{2}, \frac{1}{2}, 0, \ldots, 0)$. The minimum value of $F^*$ is strictly larger than 4 except when $u_1 = -u_2$ (equivalently, when the original signal coefficients $c_1, c_2$ are orthogonal).

2. Suppose the points $u_1, \ldots, u_N$ do not lie on a semi-circle. Then their convex hull contains the origin in the interior and in particular $N > 2$. We can find three of these $N$ points, say $u_1, u_2, u_3$, such that the triangle connecting these three points contains the origin. Then, $F^*$ is minimized at $\lambda = (\lambda_1, \lambda_2, \lambda_3, 0, \ldots, 0)$, where $\lambda_1, \lambda_2, \lambda_3$ are

---

[31]$\phi_i \in [\pi, 2\pi)$ can be replaced by $\phi_i - \pi$, corresponding to replacing the vector $c_i$ by $-c_i$.

unique weights such that $\lambda_1 u_1 + \lambda_2 u_2 + \lambda_3 u_3 = \mathbf{0}$. In this case the minimum value of $F^*$ is exactly 4.

We note that in the latter case, whenever $N > 3$, there is not a unique set of three points whose convex hull contains the origin. Thus $F^*$ is not uniquely minimized, and we cannot use the analogue of Lemma 4 to characterize society's $t$-optimal divisions.