## **Sequential Cooperative Bayesian Inference**

Junqi Wang 1 Pei Wang 1 Patrick Shafto 1

## **Abstract**

Cooperation is often implicitly assumed when learning from other agents. Cooperation implies that the agent selecting the data, and the agent learning from the data, have the same goal, that the learner infer the intended hypothesis. Recent models in human and machine learning have demonstrated the possibility of cooperation. We seek foundational theoretical results for cooperative inference by Bayesian agents through sequential data. We develop novel approaches analyzing consistency, rate of convergence and stability of Sequential Cooperative Bayesian Inference (SCBI). Our analysis of the effectiveness, sample efficiency and robustness show that cooperation is not only possible in specific instances but theoretically well-founded in general. We discuss implications for human-human and human-machine cooperation.

## 1. Introduction

Learning often occurs sequentially, as opposed to in batch, and from data provided by other agents, as opposed to from a fixed random sampling process. The canonical example of sequential learning from an agent occurs in educational contexts where the other agent is a teacher whose goal is to help the learner. However, instances appear across a wide range of contexts including informal learning, language, and robotics. In contrast with typical contexts considered in machine learning, it is reasonable to expect the cooperative agent to adapt their sampling process after each trial, consistent with the goal of helping the learner learn more quickly. It is also reasonable to expect that learners, in dealing with such cooperative agents, would know the other agent intends to cooperate and incorporate that knowledge when updating their beliefs. In this paper, we analyze basic statistical properties of such sequential cooperative inferences.

Proceedings of the 37<sup>th</sup> International Conference on Machine Learning, Vienna, Austria, PMLR 119, 2020. Copyright 2020 by the author(s).

Large behavioral and computational literatures highlight the importance cooperation for learning. Across behavioral sciences, cooperative information sharing is believed to be a core feature of human cognition. Education, where a teacher selects examples for a learner, is perhaps the most obvious case. Other examples appear in linguistic pragmatics (Frank & Goodman, 2012), in speech directed to infants (Eaves Jr et al., 2016), and children's learning from demonstrations (Bonawitz et al., 2011). Indeed, theorists have posited that the ability to select data for and learn cooperatively from others explains humans' ability to learning quickly in childhood and accumulate knowledge over generations (Tomasello, 1999; Csibra & Gergely, 2009).

Across computational literatures, cooperative information sharing is also believed to be central to human-machine interaction. Examples include pedagogic-pragmatic value alignment in robotics (Fisac et al., 2017), cooperative inverse reinforcement learning (Hadfield-Menell et al., 2016), machine teaching (Zhu, 2013), and Bayesian teaching (Eaves Jr et al., 2016) in machine learning, and Teaching dimension in learning theory (Zilles et al., 2008; Doliwa et al., 2014). Indeed, rather than building in knowledge or training on massive amounts of data, cooperative learning from humans is a strong candidate for advancing machine learning theory and improving human-machine teaming more generally.

While behavioral and computational research makes clear the importance of cooperation for learning, we lack mathematical results that would establish statistical soundness. In the development of probability theory, proofs of consistency and rate of convergence were celebrated results that put Bayesian inference on strong mathematical footing (Doob, 1949). Moreover, establishment of stability with respect to mis-specification ensured that theoretical results could apply despite the small differences between the model and reality (Kadane et al., 1978; Berger et al., 1994). Proofs of consistency, convergence, and stability ensure that intuitions regarding probabilistic inference were formalized in ways that satisfied basic desiderata.

Our goal is to provide a comparable foundation for sequential Cooperative Bayesian Inference as statistical inference for understanding the strengths, limitations, and behavior of cooperating agents. Grounded strongly in machine learning (Murphy, 2012; Ghahramani, 2015) and human learning

<sup>&</sup>lt;sup>1</sup>CoDaS Lab, Department of Math & CS, Rutgers University at Newark, New Jersey, USA. Correspondence to: Junqi Wang <junqi.wang@rutgers.edu>.

(Tenenbaum et al., 2011), we adopt a probabilistic approach. We approach consistency, convergence, and stability using a combination of new analytical and empirical methods. The result will be a model agnostic understanding of whether and under what conditions sequential cooperative interactions result in effective and efficient learning.

Notations are introduced at the end of this section. Section 2 introduces the model of sequential cooperative Bayesian inference (SCBI), and Bayesian inference (BI) as the comparison. Section 3 presents a new analysis approach which we apply to understanding consistency of SCBI. Section 4 presents empirical results analyzing the sample efficiency of SCBI versus BI, showing convergence of SCBI is considerably faster. Section 5 presents the empirical results testing robustness of SCBI to perturbations. Section 6 introduces an application of SCBI in Grid world model. Section 7 describes our contributions in the context of related work, and Section 8 discusses implications for machine learning and human learning.

**Preliminaries.** Throughout this paper, for a vector  $\theta$ , we denote its *i*-th entry by  $\theta_i$  or  $\theta(i)$ . Similarly, for a matrix M, we denote the vector of i-th row by  $M_{(i)}$ , the vector of j-th column by  $\mathbf{M}_{(-,j)}$ , and the entry of i-th row and j-th column by  $\mathbf{M}_{(i,j)}$  or simply  $\mathbf{M}_{ij}$ . Further, let  $\mathbf{r}, \mathbf{c}$  be the column vectors representing the row and column marginals (sums along row/column) of M. Let  $e_n$  or simply e be the vector of ones. The symbol  $\mathcal{N}_{\text{vec}}\left(\theta,s\right)$  is used to denote the normalization of a non-negative vector  $\theta$ , i.e.,  $\mathcal{N}_{\text{vec}}(\theta, s) =$  $\frac{s}{\sum \theta_i} \theta$  with s=1 if absent. Similarly, the normalization of matrices are denoted by  $\mathcal{N}_{col}(\mathbf{M}, \theta)$ , with "col" indicating column normalization (for row normalization, write "row" instead), and  $\theta$  denotes to which vector of sums the matrix is normalized. The set of probability distributions on a finite set  $\mathcal{X}$  is denoted by  $\mathcal{P}(\mathcal{X})$ , we do not distinguish it with the simplex  $\Delta^{|\mathcal{X}|-1}$ . The language of statistical models and estimators follows the notations of the book (Miescke & Liese, 2008).

#### 2. The Construction

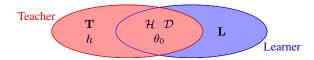


Figure 1. Two agents and their knowledge before starting.

In this paper, we consider cooperative communication models with two agents, which we call a teacher and a learner. Let  $\mathcal{H} = \{1, \dots, m\}$  be the set of m hypotheses, i.e., concepts to teach. The shared **goal** is for the learner to infer the correct hypothesis  $h \in \mathcal{H}$  which is only known by the teacher at the beginning. To facilitate learning, the teacher

passes one element from a finite data set  $\mathcal{D} = \{1, \dots, n\}$  sequentially. Each agent has knowledge about the relation between  $\mathcal{H}$  and  $\mathcal{D}$ , in terms of a positive matrix whose normalization can be treated as the likelihood matrix in a Bayesian sense. Let  $\mathbf{T}, \mathbf{L} \in \mathrm{Mat}_{n \times m}(\mathbb{R}^+)$  be the matrices for teacher and learner, respectively.

In order to construct a Bayesian theory, the learner has an initial prior  $\theta_0$  on  $\mathcal{H}$ , which, along with posteriors  $\theta_k(k\geq 1)$ , are elements in  $\mathcal{P}(\mathcal{H})=\Delta^{m-1}:=\{\theta\in\mathbb{R}^m:\sum_{i=1}^m\theta(i)=1\}$ . Privately, the teacher knows the true hypothesis  $h\in\mathcal{H}$  to teach. To measure how well a posterior  $\theta_k$  performs, we may view h as a distribution on  $\mathcal{H}$ , namely  $\widehat{\theta}=\delta_h\in\mathcal{P}(\mathcal{H})$ , and calculate the  $L^1$ -distance  $||\theta_k-\widehat{\theta}||_1$  on  $\mathcal{P}(\mathcal{H})=\Delta^{m-1}\subseteq\mathbb{R}^m$ .

We assume that  $\mathcal{H}$ ,  $\mathcal{D}$ ,  $\mathbf{T}$ ,  $\mathbf{L}$  and  $\theta_0$  satisfy:

- (i) There are no fewer data than hypotheses  $(n \ge m)$ .
- (ii) The hypotheses are distinguishable, i.e., there is no  $\lambda \in \mathbb{R}$  such that  $\mathbf{T}_{(...i)} = \lambda \mathbf{T}_{(...i)}$  for any  $i \neq j$ , and so is  $\mathbf{L}$ .
- (iii) T is a *scaled matrix* of L, i.e., there exist invertible diagonal matrices  $E_1$  and  $E_2$  such that  $T = E_1 L E_2$ . (Both agents aware this assumption, though possibly neither know the other's matrix.)
- (iv)  $\theta_0$  is known by the teacher.

Our model is constructed and studied under these assumptions (Sec. 3 and Sec. 4). We also studied stability under violations of (iii) and (iv), where we assume that  $\mathbf{T}$  and teacher's knowledge  $\theta_0^T$  about  $\theta_0$  is slightly different from (some scaled matrix of)  $\mathbf{L}$  and  $\theta_0$  (Sec. 5). Assumption (iii) is a relaxation of the assumption of Bayesian inference that  $\mathbf{T} = \mathbf{L} = \mathbf{M}$  is the likelihood matrix. Practically, we may achieve (iii) by adding to the common ground a shared matrix  $\mathbf{M}$  (e.g. joint distribution on  $\mathcal{D}$  and  $\mathcal{H}$ ) and scaling it to  $\mathbf{T}$  and  $\mathbf{L}$ . We may obtain  $\mathbf{M}$  by taking the same ground model or using the same statistical data (e.g. historical statistical records). In fact, with (iii), it does not affect the process of inference whether  $\mathbf{M}$  is accessible to agents.

In SCBI (see details in later this section), thanks to the property that a matrix and its scaled matrices behave the same in Sinkhorn scaling (Hershkowitz et al., 1988), the pre-processings of **T** and of **L** lead to the same results under (iii) and (iv). Thus assumption (iii) is equivalent to:

(iii') T = L = M where M is a column-stochastic matrix.

We assume (iii') is valid until we discuss stability.

In our setup, the teacher teaches in sequence. At each round the teacher chooses a data point from  $\mathcal D$  by sampling according to a distribution. And the learner learns by maintaining a posterior distribution on  $\mathcal H$  through Bayesian inference with likelihood matrices not necessarily fixed.

Formally, the teacher's job is to select a sequence of data  $(d_k)_{k\in\mathbb{N}}$  by first constructing a sequence of random variables  $(D_k)_{k\in\mathbb{N}}$ , then sampling each  $d_k$  as a realization of  $D_k$ . Each  $d_k$  is given to the learner at round k. And the learner's job is to produce a sequence of posteriors  $(\theta_k)_{k\in\mathbb{N}}$  on  $\mathcal{P}(\mathcal{H})$ . To calculate  $\theta_k$ , learner can use the matrix  $\mathbf{L}$ , the initial prior  $\theta_0$  which is common knowledge, and the sequence of data  $(d_i)_{i\leq k}$  which is visible at round k. The learner find each posterior by giving a function  $S_k((d_i)_{i\leq k};\mathbf{L},\theta_0)^{-1}$  for k>0. We may further define  $S_0(\varnothing;\mathbf{L},\theta_0)=\theta_0$ .

Since  $(d_k)_{k\in\mathbb{N}}$  is generated by a sequence of random variables  $(D_k)_{k\in\mathbb{N}}$ , the function  $S_k$  can be treated as a function taking  $(D_i)_{i\leq k}$  as inputs and producing a random variable  $\Theta_k$  as output. We call the distribution of  $\Theta_k$  by  $\mu_k\in\mathcal{P}(\mathcal{P}(\mathcal{H}))=\mathcal{P}(\Delta^{m-1})$ . The  $S_k$ 's as functions of random variables are called *estimators*.

Being a special case of the above framework, Bayesian inference dealing with sequential data is a well-studied model. However, there is no cooperation in Bayesian inference since the teaching distribution and learning likelihood are constant on time (the teacher side is typically left implicit). To introduce cooperation following cooperative inference (Yang et al., 2018), we propose Sequential Cooperative Bayesian Inference (SCBI), which is a sequential version of the cooperative inference.

### 2.1. Sequential Cooperative Bayesian Inference

Sequential Cooperative Bayesian Inference (SCBI) assumes that the two agents—a teacher and a learner—cooperate to facilitate learning. Prior research has formalized this cooperation (in a single-round game) as a system of two interrelated equations in which the teacher's choice of data depends on the learner's inference, and the learner's inference depends on reasoning about the teacher's choice. This prior research into such Cooperative Inference has focused on batch selection of data (Yang et al., 2018; Wang et al., 2019a), and has been shown to be formally equivalent to Sinkhorn scaling (Wang et al., 2019b). Following this principle, we propose a new *sequential* setting in which the teacher chooses data sequentially, and both agents update the likelihood at each round to optimize learning.

**Cooperative Inference.** Let  $P_{L_0}(h)$  be the learner's prior of hypothesis  $h \in \mathcal{H}$ ,  $P_{T_0}(d)$  be the teacher's prior of selecting data  $d \in \mathcal{D}$ . Let  $P_T(d|h)$  be the teacher's likelihood of selecting d to convey h and  $P_L(h|d)$  be the learner's posterior for h given d. **Cooperative inference** is then a system of two equations shown below, with  $P_L(d)$  and  $P_T(h)$  the normalizing constants:

$$P_L(h|d) = \frac{P_T(d|h) P_{L_0}(h)}{P_L(d)}, \ P_T(d|h) = \frac{P_L(h|d) P_{T_0}(d)}{P_T(h)}. \ (1)$$

```
Algorithm 1 SCBI, without assumption (iii')

== Teacher's Part: ==
Input: \mathbf{T} \in \operatorname{Mat}_{n \times m}(\mathbb{R}^+), \theta_0, h \in \mathcal{H}, (\widehat{\theta} = \delta_h)

Output: Share (d_1, d_2, \dots) to learner
for all i \geq 1 do

sample d_i \sim \mathscr{N}_{\text{vec}}\left(\mathbf{T}^{\langle n\theta_{i-1} \rangle}_{(-,h)}, 1\right).

\theta_i \leftarrow \mathbf{T}^{\langle n\theta_{i-1} \rangle}_{(d_i,-)} estimation of learner's posterior end for

== Learner's Part: ==
Input: \mathbf{L} \in \operatorname{Mat}_{n \times m}(\mathbb{R}^+), \theta_0, (d_1, d_2, \dots)

Output: (\theta_0, \theta_1, \theta_2, \dots) posteriors
for all i \geq 1 do

\theta_i \leftarrow \mathbf{L}^{\langle n\theta_{i-1} \rangle}_{(d_i,-)}
end for
```

Note:  $\mathbf{T}^{\langle n\theta_{i-1} \rangle}$ ,  $\mathbf{L}^{\langle n\theta_{i-1} \rangle}$  are the  $\mathbf{M}^{\langle \mathbf{c}_{k-1} \rangle}$  in the text.

It is shown (Wang et al., 2019a;b) that Eq. (1) can be solved using Sinkhorn scaling, where  $(\mathbf{r}, \mathbf{c})$ -Sinkhorn scaling of a matrix  $\mathbf{M}$  is simply the iterated alternation of row normalization of  $\mathbf{M}$  with respect to  $\mathbf{r}$  and column normalization of  $\mathbf{M}$  with respect to  $\mathbf{c}$ . The limit of such iterations exist if the sums of elements in  $\mathbf{r}$  and  $\mathbf{c}$  are the same (Schneider, 1989).

**Sequential Cooperation.** SCBI allows multiple rounds of teaching and requires each choice of data to be generated based on cooperative inference, with the learner updating their beliefs between each round. In each round, based on the data being taught and the learner's initial prior on  $\mathcal{H}$  as common knowledge, the teacher and learner update their common likelihood matrix according to cooperative inference (using Sinkhorn scaling), then the data selection and inference proceed based on the updated likelihood matrix.

Precisely, starting from learner's prior  $S_0 = \theta_0 \in \Delta^{m-1}$ , let the data been taught up to round k be  $(d_1, \ldots, d_{k-1})$  and the posterior of the learner after round k-1 be  $\theta_{k-1} =$  $S_{k-1}(d_1,\ldots,d_{k-1};\theta_0) \in \mathcal{P}(\mathcal{H})$ , which is actually predictable for both agents (obvious for k = 1 and inductively correct for k > 1 by later argument). To teach, the teacher calculates the Sinkhorn scaling of M given the uniform row sums  $\mathbf{r}_{k-1} = \mathbf{e}_n = (1, 1, \dots, 1)^{\top}$  and column sums  $\mathbf{c}_{k-1} = n\theta_{k-1}$  (to make the sum of  $\mathbf{r}_{k-1}$  equal that of  $c_{k-1}$ , which guarantees the existence of the limit in Sinkhorn scaling), denoted by  $\mathbf{M}^{\langle \mathbf{c}_{k-1} \rangle}$ . The teacher's data selection is proportional to columns of  $\mathbf{M}^{\langle \mathbf{c}_{k-1} \rangle}$ . Thus let  $\mathbf{M}_k$  be the column normalization of  $\mathbf{M}^{\langle \mathbf{c}_{k-1} \rangle}$  by  $\mathbf{e}_m$ , i.e.,  $\mathbf{M}_k = \mathscr{N}_{\text{col}}(\mathbf{M}^{\langle \mathbf{c}_{k-1} \rangle}, \mathbf{e}_m)$ . Then the teacher defines  $\mathbf{D}_k$ using distribution  $(\mathbf{M}_k)_{(.,h)}$  on set  $\mathcal{D}$  and samples  $d_k \sim D_k$ , then passes  $d_k$  to the learner.

On learner's side, the learner obtains the likelihood matrix  $\mathbf{M}_k$  in the same way as above and applies normal Bayesian inference with datum  $d_k$  past from the teacher. First, learner

<sup>&</sup>lt;sup>1</sup>we may omit **L** and (or)  $\theta_0$  when there is no ambiguity.

takes the prior to be the posterior of the last round,  $\theta_{k-1} = \frac{1}{n}\mathbf{c}_{k-1}$ , then multiply it by the likelihood of selecting  $d_k$ —the row of  $\mathbf{M}_k$  corresponding to  $d_k$ , which results  $\mathring{\theta}_k = (\mathbf{M}_k)_{(d_k,\cdot)} \mathrm{diag}(\theta_{k-1})$ . Then the posterior  $\theta_k$  is obtained by row normalizing  $\mathring{\theta}_k$ . Inductively, in the next round, the learner will start with  $\theta_k$  and  $\mathbf{c}_k = n\theta_k$ . The learner's calculation in round k can be simulated by the teacher, so the teacher can predict  $\theta_k$ , which inductively shows the assumption (teacher knows  $\theta_{k-1}$ ) in previous paragraph.

The calculation can be simplified. Consider that the vector  $\mathbf{c}_{k-1}$ , being proportional to the prior, is used in  $\mathbf{M}_k = \mathcal{N}_{\mathrm{col}}(\mathbf{M}^{\langle \mathbf{c}_{k-1} \rangle}, \mathbf{e}_m) = \mathbf{M}^{\langle \mathbf{c}_{k-1} \rangle} \left( \mathrm{diag}(n\theta_{k-1}) \right)^{-1}$ , then  $\mathring{\theta}_k = \left( \mathbf{M}^{\langle \mathbf{c}_{k-1} \rangle} \left( \mathrm{diag}(n\theta_{k-1}) \right)^{-1} \mathrm{diag}(\theta_{k-1}) \right)_{(d_k,.)} = \frac{1}{n} \mathbf{M}^{\langle \mathbf{c}_{k-1} \rangle}$ . Furthermore, since  $\mathbf{M}^{\langle \mathbf{c}_{k-1} \rangle}$  is row normalized to  $\mathbf{e}_m$ , each row of it is a probability distribution on  $\mathcal{H}$ . Thus  $S_k(d_1,\ldots,d_k) = \theta_k = n\mathring{\theta}_{k-1} = \mathbf{M}^{\langle \mathbf{c}_k \rangle} \left( \frac{1}{d_k} \right)^{2}$ .

The simplified version of SCBI algorithm is given in Algorithm 1.

#### 2.2. Bayesian Inference: the Control

In order to test the performance of SCBI, we recall the classical Bayesian inference (BI). In BI, a fixed likelihood matrix  $\mathbf{M}$  is used throughout the communication process. Bayes' rule requires  $\mathbf{M}$  to be the conditional distribution on the set of data given each hypothesis, thus  $\mathbf{M} = \mathbf{T} = \mathbf{L}$  is column-stochastic as in assumption (iii').

For the teacher, given  $h \in \mathcal{H}$ , the teaching distribution is the column vector  $P_h = \mathbf{M}_{(\cdot,h)} \in \mathcal{P}(\mathcal{D})$ . This defines random variable  $\mathbf{D}_k$ . Then the teacher selects data via i.i.d. sampling according to  $P_h$ . The random variables  $(\mathbf{D}_k)_{k\geq 1}$  are identical.

The learner first chooses a prior  $\theta_0 \in \mathcal{P}(\mathcal{H})$  ( $\theta_0 = S_0$  is part of the model, usually the uniform distribution), then uses Bayes' rule with likelihood  $\mathbf{M}$  to update the posterior distribution repeatedly. Given taught datum d, the map from the prior  $\theta$  to the posterior distribution is denoted by  $B_d(\theta) = \mathcal{N}_{\text{vec}}\left(\mathbf{M}_{(d,.)}\mathrm{diag}(\theta),1\right)$ . Thus the learner's estimation over  $\mathcal{H}$  given a sequential data  $(d_1,\ldots,d_k)$  can be written recursively by  $S_0=\theta_0$ , and  $S_k(d_1,\ldots,d_k)=B_{d_k}(S_{k-1}(d_1,\ldots,d_{k-1}))$ . Thus, by induction,  $S_k(d_1,\ldots,d_k)=(B_{d_k}\circ B_{d_{k-1}}\circ\cdots\circ B_{d_1})(S_0)$ .

#### 3. Consistency

We investigate the effectiveness of the estimators in both BI and SCBI by testing their *consistency*: setting the true hypothesis  $h \in \mathcal{H}$ , given  $(D_k)$ ,  $(S_k)$  and  $\theta_0$ , we examine the convergence (using the  $L^1$ -distance on  $\mathcal{P}(\mathcal{H})$ ) of the

posterior sequence  $(\Theta_k) = (S_k(D_1, \dots, D_k))$  as sequence of random variables and check whether the limit is  $\widehat{\theta}$  as a constant random variable.

#### 3.1. BI and KL Divergence

The consistency of BI has been well studied since Bernstein and von Mises and Doob (Doob, 1949). In this section, we state it in our situation and derive a formula for the rate of convergence, as a baseline for the cooperative theory. Derivations and proofs can be found in the Supplementary Material.

**Theorem 3.1.** [(Miescke & Liese, 2008, Theorem 7.115)] In BI, the sequence of posteriors  $(S_k)$  is strongly consistent at  $\hat{\theta} = \delta_h$  for each  $h \in \mathcal{H}$ , with arbitrary choice of an interior point  $\theta_0 \in (\mathcal{P}(\mathcal{H}))^{\circ}$  (i.e.  $\theta_0(h) > 0$  for all  $h \in \mathcal{H}$ ) as prior.

Remark 1. For a fixed true distribution  $\widehat{\theta}$ , strong consistency of  $(S_k)_{k\in\mathbb{N}}$  is defined to be: the sequence of posteriors  $\Theta_k$  given by the estimator  $S_k$ , as a sequence of random variables, converges to  $\widehat{\theta}$  (as a constant random variable) almost surely according to random variables  $(D_k)_{k\in\mathbb{N}}$  that the teacher samples from. If the convergence is in probability, the sequence of estimators is said to be *consistent*.

Remark 2. Theorem 3.1 also assumes that hypotheses are distinguishable (Section 2). In a general theory of statistical models,  $\widehat{\theta}$  is not necessarily  $\delta_h$  for some  $h \in \mathcal{H}$ . However, in BI, it is critical to have  $\widehat{\theta} = \delta_h$ , since BI with a general  $\widehat{\theta} \in \mathcal{P}(\mathcal{H})$  is almost never consistent or strongly consistent.

Consistency—independent of the choice of prior  $\theta_0$  interior of  $\mathcal{P}(\mathcal{H})$ —guarantees that BI is always effective.

**Rate of Convergence.** After effectiveness, we provide the efficiency of BI in terms of asymptotic rate of convergence.

**Theorem 3.2.** In BI, with  $\widehat{\theta} = \delta_h$  for some  $h \in \mathcal{H}$ , let  $\Theta_k(h)(D_1,\ldots,D_k) := S_k(h|D_1,\ldots,D_k)$  be the h-component of posterior given  $D_1,\ldots,D_k$  as random variables valued in  $\mathcal{D}$ . Then  $\frac{1}{k}\log\left(\frac{\Theta_k(h)}{1-\Theta_k(h)}\right)$  converges to a constant  $\min_{h'\neq h}\left\{\mathrm{KL}(\mathbf{M}_{(.,h)},\mathbf{M}_{(.,h')})\right\}$  almost surely. Remark 3. We call  $\min_{h'\neq h}\left\{\mathrm{KL}(\mathbf{M}_{(.,h)},\mathbf{M}_{(.,h')})\right\}$  the asymptotic rate of convergence (RoC) of BI, denoted by  $\mathfrak{R}^b(\mathbf{M};h)$ .

## 3.2. SCBI as a Markov Chain

From the proof of Theorem 3.1, the pivotal property is that the variables  $D_1, D_2, \ldots$  are commutative in posteriors (the variables can occur in any order without affecting the posterior) thanks to commutativity of multiplication. However, in SCBI, the commutativity does not hold, since the likelihood matrix depends on previous outcome. Thus the method used in BI analysis no longer works here.

<sup>&</sup>lt;sup>2</sup>See Supplementary Material for detailed examples.

Because the likelihood matrix  $\mathbf{M}_k = \mathbf{M}^{\langle \mathbf{c}_{k-1} \rangle}$  depends on the predecessive state only, the process is in fact Markov, we may analyze the model as a Markov chain on the continuous state space  $\mathcal{P}(\mathcal{H})$ .

To describe this process, let  $\mathcal{P}(\mathcal{H}) = \Delta^{m-1}$  be the space of states, and let  $h \in \mathcal{H}$  be the true hypothesis to teach  $(\widehat{\theta} = \delta_h)$ , let learner's prior be  $S_0 = \theta_0 \in \mathcal{P}(\mathcal{H})$ , or say, the distribution of learner's initial state is  $\mu_0 = \delta_{\theta_0} \in \mathcal{P}(\mathcal{P}(\mathcal{H}))$ .

**The operator**  $\Psi$ **.** In the Markov chain, in each round, the transition operator maps the prior as a probability distribution on state space  $\mathcal{P}(\mathcal{H}) = \Delta^{m-1}$  to the posterior as another, i.e.,  $\Psi(h): \mathcal{P}(\mathcal{P}(\mathcal{H})) \to \mathcal{P}(\mathcal{P}(\mathcal{H}))$ .

To make the formal definition of  $\Psi(h)$  simpler, we need to define some maps. For any  $d \in \mathcal{D}$ , let  $T_d : \Delta^{m-1} \to \Delta^{m-1}$ be the map bringing the learner's prior to posterior when data d is chosen by the teacher, that is,  $T_d$  sends each normalized vector  $\theta$  to  $T_d(\theta) = \mathbf{M}^{\langle n\theta \rangle}_{(d,.)}$  according to SCBI. Each  $T_d$  is a bijection based on the uniqueness of Sinkhorn scaling limits of M, shown in (Hershkowitz et al., 1988). Further, the map  $T_d$  is continuous on  $\Delta^{m-1}$  and smooth in its interior according to (Wang et al., 2019b). Continuity and smoothness of  $T_d$  make it natural to induce a push-forward  $T_{d*}: \mathcal{P}(\Delta^{m-1}) \to \mathcal{P}(\Delta^{m-1})$  on Borel measures. Explicitly,  $(T_{d*}(\mu))(E) = \mu(T_d^{-1}(E))$  for each Borel measure  $\mu \in \mathcal{P}(\Delta^{m-1})$  and each Borel measurable set  $E \subseteq \Delta^{m-1}$ . Let  $\tau : \mathcal{P}(\mathcal{H}) \to \mathcal{P}(\mathcal{D})$  be the map of teacher's adjusting sample distribution based on the learner's prior, that is, given a learner's prior  $\theta \in \Delta^{m-1}$ , by definition of SCBI, the distribution of the teacher is adjusted to  $\tau(\theta) = \frac{\mathbf{M}^{(n\theta)}_{(\cdot,h)}}{n\theta(h)} = (n\theta(h))^{-1}(T_1(\theta)(h),\dots,T_n(\theta)(h)).$  Each component d of  $\tau$  is denoted by  $\tau_d$ . We can use  $\tau$  only for  $\theta_0 = \delta_h$  in which case teacher can trace learner's state. Now we can define  $\Psi(h)$  formally.

**Definition 3.3.** Given a fixed hypothesis  $h \in \mathcal{H}$ , or say  $\delta_h \in \mathcal{P}(\mathcal{H})$ , the operator  $\Psi(h) : \mathcal{P}(\Delta^{m-1}) \to \mathcal{P}(\Delta^{m-1})$  translating a prior as a Borel measure  $\mu$  to the posterior distribution  $\Psi(h)(\mu)$  according to one round of SCBI is given below, for any Borel measurable set  $E \subset \Delta^{m-1}$ .

$$(\Psi(h)(\mu))(E) := \int_{E} \sum_{d \in \mathcal{D}} \tau_d(T_d^{-1}(\theta)) d(T_{d*}(\mu))(\theta).$$
 (2)

In our case, we start with a distribution  $\delta_{\theta}$  where  $\theta \in \mathcal{P}(\mathcal{H})$  is the prior of the learner on the set of hypotheses. In each round of inference, there are n different possibilities according to the data taught. Thus in any finite round k, the distribution of the posterior is the sum of at most  $n^k$  atoms (actually, we can prove  $n^k$  is exact). Thus in the following discussions, we assume that  $\mu$  is atomic. The  $\Psi$  action on an atomic distribution is determined by that of an atom:

$$\Psi(h)(\delta_{\theta}) = \sum_{i=1}^{n} \frac{\mathbf{M}_{(i,h)}^{(n\theta)}}{n\theta(h)} \delta_{(\mathbf{M}_{(i,-)}^{(n\theta)})}.$$
 (3)

Moreover, since the SCBI behavior depends only on the prior (with fixed M and h) as a random variable, the same operator  $\Psi(h)$  applies to every round in SCBI. Thus we can conclude that the following proposition is valid:

**Proposition 3.4.** Given  $h \in \mathcal{H}$ , let  $\widehat{\theta} = \delta_h$ , the sequence of estimators  $(S_k)_{k \in \mathbb{N}}$  in SCBI forms a time-homogeneous Markov chain on state space  $\mathcal{P}(\mathcal{H})$  with transition operator  $\Psi(h)$  characterized by Eq. (2) and Eq. (3).

Thanks to the fact that the SCBI is a time homogeneous Markov process, we can further show the consistency.

**Theorem 3.5** (Consistency). In SCBI, let M be a positive matrix. If the teacher is teaching one hypothesis h (i.e.,  $\widehat{\theta} = \delta_h \in \mathcal{P}(\mathcal{H})$ ), and the prior distribution  $\mu_0 \in \mathcal{P}(\Delta^{m-1})$  satisfies  $\mu_0 = \delta_{\theta_0}$  with  $\theta_0(h) > 0$ , then the estimator sequence  $(S_k)$  is consistent, for each  $h \in \mathcal{H}$ , i.e., the posterior random variables  $(\Theta_k)_{k \in \mathbb{N}}$  converge to the constant random variable  $\widehat{\theta}$  in probability.

Remark 4. The assumption in Theorem 3.5 that  $\theta_0(h) > 0$  is necessary in any type of Bayesian inference since it is impossible to get the correct answer in posterior by Bayes' rule, if it is excluded in the prior at the beginning. In practice, the prior distribution is usually chosen to be  $\mu_0 = \delta_{\mathbf{u}}$  with the uniform distribution vector in  $\mathcal{P}(\mathcal{H})$ , i.e.,  $\mathbf{u} = \frac{1}{m}(1,\dots,1)^{\top} \in \Delta^{m-1}$ .

**Rate of Convergence.** Thanks to consistency, we can calculate the asymptotic rate of convergence for SCBI.

**Theorem 3.6.** With matrix M, hypothesis  $h \in \mathcal{H}$ , and a prior  $\mu_0 = \delta_{\theta_0} \in \mathcal{P}(\Delta^{m-1})$  same as in Theorem. 3.5, let  $\theta_k$  denote a sample value of the posterior  $\Theta_k$  after k rounds of SCBI, then

$$\lim_{k \to \infty} \mathbb{E}_{\mu_k} \left[ \frac{1}{k} \log \left( \frac{\theta_k(h)}{1 - \theta_k(h)} \right) \right] = \mathfrak{R}^{\mathrm{s}}(\mathbf{M}; h) \tag{4}$$

where  $\mathfrak{R}^{s}(\mathbf{M};h) := \min_{h \neq h'} \mathrm{KL}\left(\mathbf{M}_{(.,h)}^{\sharp}, \mathbf{M}_{(.,h')}^{\sharp}\right)$  with  $\mathbf{M}^{\sharp} = \mathscr{N}_{col}(\mathrm{diag}(\mathbf{M}_{(.,h)})^{-1}\mathbf{M})$ . Thus we call  $\mathfrak{R}^{s}(\mathbf{M};h)$  the asymptotic rate of convergence (RoC) of SCBI.

## 4. Sample Efficiency

In this section, we present some empirical results comparing the sample efficiency of SCBI and BI.

#### 4.1. Asymptotic RoC Comparison

We first compare the asymptotic rate of convergence ( $\mathfrak{R}^{b}$  for BI and  $\mathfrak{R}^{s}$  for SCBI, see Theorems 3.2 and 3.6). The matrix **M** is sampled through m i.i.d. uniform distributions on  $\Delta^{n-1}$ , one for each column.

For each column-normalized matrix  $\mathbf{M}$ , we compute two variables to compare BI with SCBI: the probability  $\mathfrak{P} := \Pr\left(\frac{1}{m} \sum_{h \in \mathcal{H}} \mathfrak{R}^{\mathrm{s}}(\mathbf{M}; h) \geq \frac{1}{m} \sum_{h \in \mathcal{H}} \mathfrak{R}^{\mathrm{b}}(\mathbf{M}; h)\right)$ 

and the expected value of averaged difference  $\mathfrak{E} := \mathbb{E} \left[ \frac{1}{m} \sum_{h \in \mathcal{H}} \mathfrak{R}^{\mathrm{s}}(\mathbf{M}; h) - \frac{1}{m} \sum_{h \in \mathcal{H}} \mathfrak{R}^{\mathrm{b}}(\mathbf{M}; h) \right].$ 

**Two-column Cases.** Consider the case where  $\mathbf{M}$  is of shape  $n \times 2$  with the two columns sampled from  $\Delta^{n-1}$  uniformly and independently, we simulated for  $n=2,3,\ldots,50$  with a size- $10^{10}$  Monte Carlo method for each n to calculate  $\mathfrak{P}$  and  $\mathfrak{E}$ . The result is shown in Fig. 2(A)(B).

We can reduce the calculation of  $\mathfrak{E}$  to a numerical integral  $\mathfrak{E} = \int_{(\Delta^{n-1})^2} \ln \left( \sum_{i=1}^n \frac{\mathbf{x}_i}{\mathbf{y}_i} \right) d\mathbf{x} d\mathbf{y} - \ln n - \frac{n-1}{n}.$ 

Since  $\mathfrak{P}$  goes too close to 1 as the rank grows, we use  $-\ln(1-\mathfrak{P})$  to show the increasing in detail. <sup>4</sup>

More Columns of a Fixed Row Size. To verify the general cases, we simulated  $\mathfrak{P}$  and  $\mathfrak{E}$  by Monte Carlo on matrices of 10-row and various-column shapes, see Fig. 2(C)(D). We sampled  $10^8$  different M of shape  $10 \times m$  for each  $2 \leq m \leq 10$ . Empirical results show that  $\mathfrak{E}$  decreases slowly but  $\mathfrak{P}$  still increase logistically as m grows.

**Square Matrices.** Fig. 2(E)(F) shows the square cases with size from 2 to 50, simulated by size  $10^8$  Monte Carlo.

The empirical  $\mathfrak P$  is the mean of N (sample-size) i.i.d. variables valued 0 or 1, thus the standard deviation of a single variable is smaller than 1. By Central Limit Theorem, the standard deviation  $\sigma(\mathfrak P) < N^{-1/2}$  (precision threshold). So we draw lines  $y = N^{-1/2}$  in each log-figure, but only in one figure the line lies in the view area.

In all simulated cases, we observe that  $\mathfrak{E}>0$  and  $\mathfrak{P}>0.5$ , indicating that SCBI converges faster than BI in most cases and in average. It is also observed that SCBI behaves even better as matrix size grows, especially when the teacher has more choices on the data to be chosen (i.e., more rows).

#### 4.2. Distribution of Inference Results

The promises of cooperation is that one may infer hypotheses from small amounts of data. Hence, we compare SCBI with BI after small, fixed numbers of rounds.

We sample matrices of shape  $20 \times 20$  whose columns are distributed evenly in  $\Delta^{19}$  to demonstrate. Equivalently, they are column-normalizations of the uniformly sampled matrices whose sum of all entries is one.

Assume that the correct hypothesis to teach is  $h \in \mathcal{P}(\mathcal{H})$  We first simulate a 5-round inference behavior, exploring all possible ways that the teacher may teach, then calculate the expectation and standard deviation of  $\theta(h)$ . With 300 matrices sampled in the above way, Fig. 3 shows this comparison between BI and SCBI.

Similarly, we extend the number of rounds to 30 by Monte Carlo since an exact calculation on exhausting all possible teaching paths becomes impossible. With sampling 500 matrices independently, we simulate a teacher teaches 2000 times to round 30 for each matrix, and the statistics are also shown in Fig. 3. From Fig. 3, we observe that SCBI have better expectation and less variance in the short run.

In conclusion, experiments indicate that SCBI is both more efficient asymptotically, and in the short run.

## 5. Stability

In this section, we study the robustness of SCBI by setting the initial conditions of teacher and learner different. This could happen when agents do not have full access to their partner's exact state.

**Theory.** In this section, we no longer have assumption (iii). Let  $\mathbf{T}$  and  $\mathbf{L}$  be matrices of teacher and learner (not necessarily have (iii)). Let  $\theta_0^T$  and  $\theta_0^L$  be elements in  $\mathcal{P}(\mathcal{H})$  representing the prior on hypotheses that the teacher and learner use in the estimation of inference (teacher) and in the actual inference (learner), i.e.,  $\mu_0^T = \delta_{\theta_0^T}$  and  $\mu_0^L = \delta_{\theta_0^L}$ . During the inference, let  $\mu_k^T$  and  $\mu_k^L$  be the distribution of posteriors of the teacher and the learner after round k, and denote the corresponding random variables by  $\theta_k^T$  and  $\theta_k^L$ , for all positive k and  $\infty$ , where  $\infty$  represents the limit in probability.

Let D be a random variable on  $\mathcal{D}$ , we define an operator  $\Psi^{\mathbf{L}}_{\mathrm{D}}: \mathcal{P}(\mathcal{P}(\mathcal{H})) \longrightarrow \mathcal{P}(\mathcal{P}(\mathcal{H}))$  similar to the  $\Psi$  in Section 3. Let  $T_d(\theta) = \mathbf{L}^{\langle n\theta \rangle}_{(d, \cdot)}$ , then  $\mathrm{d}(\Psi^{\mathbf{L}}_{\mathrm{D}}(\mu))(\theta) := \sum_{d \in \mathcal{D}} \mathrm{P}(\mathrm{D} = d) \mathrm{d}(T_{d*}\mu)(\theta)$ .

**Proposition 5.1.** Given a sequence of identical independent  $\mathcal{D}$ -valued random variables  $(D_i)_{i\geq 1}$  following the uniform distribution. Let  $\mu_0\in\mathcal{P}(\mathcal{P}(\mathcal{H}))$  be a prior distribution on  $\mathcal{P}(\mathcal{H})$ , and  $\mu_{k+1}=\Psi^{\mathbf{L}}_{D_{k+1}}(\mu_k)$ , then  $\mu_k$  converges, in probability, to  $\sum_{i\in\mathcal{H}}a_i\delta_i$  where  $a_i=\mathbb{E}_{\mu_0}\left[\theta(i)\right]$ .

Remark 5. This proposition helps accelerate the simulation, that one may terminate the teaching process when  $\theta_k^T$  is sufficiently close to  $\delta_h$ , since Prop. 5.1 guarantees that the expectation of the learner's posterior on the true hypothesis h at that time is close enough to the eventual probability of getting  $\delta_h$ , i.e.  $\mathbb{E}\theta_\infty^L(h) \approx \mathbb{E}\theta_k^L(h)$ .

**Definition 5.2.** We call  $\mathbb{E}\theta_{\infty}^L(h) := \lim_{k \to \infty} \mathbb{E}_{\mu_k}(\theta(h))$  the successful rate of the inference given  $\mathbf{T}, \mathbf{L}, \theta_0^T$  and  $\theta_0^L$ . By the setup in Section 2, the failure probability,  $1 - \mathbb{E}\theta_{\infty}^L(h)$ , is  $\frac{1}{2}||\mathbb{E}\theta_{\infty}^L - \delta_h||_1$ , half of the 1-distance on  $\mathcal{P}(\mathcal{H})$ .

Simulations with Perturbation on Priors. We simulated the square cases of rank 3 and 4. We sample 5 matrices ( $\mathbf{M}_1$  to  $\mathbf{M}_5$ ) of size  $3 \times 3$ , whose columns distribute uniformly on  $\mathcal{P}(\{d_1, d_2, d_3\}) = \Delta^2$ , and 5 priors  $(\theta_1 \text{ to } \theta_5)$  in  $\mathcal{P}(\mathcal{H})$ ,

<sup>&</sup>lt;sup>3</sup>Details can be found in Supplementary Material.

<sup>&</sup>lt;sup>4</sup>We guess an empirical formula  $-\ln(1-\mathfrak{P}) \approx \frac{1}{2}\ln(x(x+1)/(x-1.5)) + 0.1x - 0.3$ , see Supplementary Material.

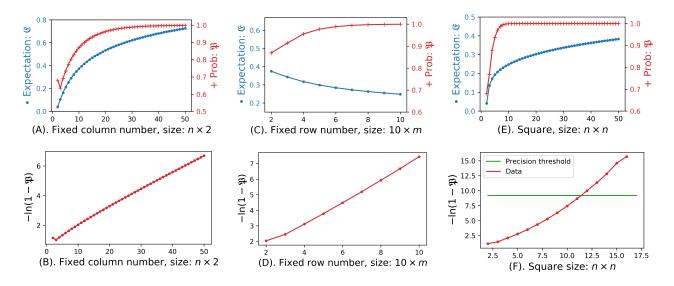


Figure 2. Comparison of RoC between BI and SCBI. (A), (C), (E): the comparison on  $\mathfrak P$  in blue and on  $\mathfrak E$  in red. (B), (D), (F): plotting  $-\ln(1-\mathfrak P)$ . (A), (B): two-column case, number of rows from 2 to 50. Monte Carlo of  $10^{10}$  samples for each point on figure. (C), (D): 10-row case, number of columns from 2 to 10. Monte Carlo of size  $10^8$ . (E), (F): square case, number of rows from 2 to 50. Monte Carlo of size  $10^8$ . The horizontal line in (F) is the theoretical threshold of precision by central limit theorem. For n > 17, MC provides  $\mathfrak P = 1$  ( $\mathfrak P^s > \mathfrak P^b$  for all samples). From the figures, except in (C) where  $\mathfrak E$  decays slowly when column number grows, the two values  $\mathfrak E$  and  $\mathfrak P$  increases as size grows in all the other cases. Moreover,  $\mathfrak P$  grows to 1 logistically in all situations.

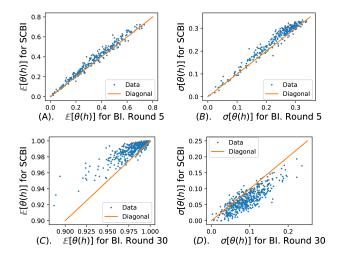


Figure 3. Comparison between BI and SCBI on  $20 \times 20$  matrices: Top: 300 points (matrices) of round 5 accurate value. Bottom: 500 points of round 30 using Monte Carlo of size 2000. Left: comparison on expectations of learner's posterior on h. Right: comparison on the standard deviations. Orange line is the diagonal.

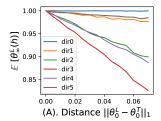
used as  $\theta_0^T$ . Similarly, we sample 3 matrices  $(\mathbf{M}_1', \mathbf{M}_2')$  and  $(\mathbf{M}_3')$  of size  $4 \times 4$ , and 3 priors  $(\theta_1', \theta_2', \theta_3')$  from  $(0.25)^3$  in the same way as above. In both cases, we assume  $(0.25)^3$  to be the true hypothesis to teach.

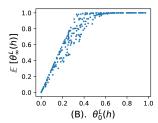
Our simulation is based on Monte Carlo method of  $10^4$  teaching sequences (for each single point plotted) then use

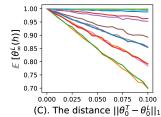
Proposition 5.1 to calculate the successful rate of inference. For  $3 \times 3$  matrices, we perturb  $\theta_0^L$  in two ways: (1) take  $\theta_0^L$ around  $\theta_0^T$  distributed evenly on concentric circles, thus 630 points for each  $\theta_0^T$  are taken. In this area, there are 84 points lying on 6 given directions (60° apart, see Supplementary Material for figures). (2) sample  $\theta_0^L$  evenly in the whole simplex  $\mathcal{P}(\mathcal{H}) = \Delta^2$  (300 points for each  $\theta_0^T$ ). For  $4 \times 4$ matrices, we perturb  $\theta_0^L$  in two ways: (1) along 15 randomly chosen directions in  $\tilde{\Delta}^3$  evenly take 21 points on each direction, and (2) sample 300 points evenly in  $\Delta^3$ . Then we have the following figure samples (for figures demonstrating the entire simulation, please see Supplementary Material). From the figures we see: 1. left pictures indicate that the learner's expected posterior on h is roughly linear to perturbations along a line. 2. right pictures indicate that the learner's expected posterior on h is closely bounded by a multiple of the learner's prior on true h. Thus we have the following conjecture:

**Conjecture 5.3.** Given  $\mathbf{L} = \mathbf{T} = \mathbf{M}$  and  $\theta_0^T$ , let h be the true hypothesis to teach. For any  $\epsilon > 0$ , let  $\theta_0^L$  be learner's prior with a distance to  $\theta_0^T$  less than  $\epsilon$ . Then the successful rate for sufficiently many rounds is greater than  $1 - C\epsilon$ , where  $C = \frac{1}{\theta_0^T(h)}$ .

Simulations with Perturbation on Matrices. We now investigate robustness of SCBI to differences between agents' matrices. Let  $\mathbf{T}$  and  $\mathbf{L}$  be stochastic, and let  $\mathbf{L}$  be perturbed from  $\mathbf{T}$ . The simulations are performed on the matrices  $\mathbf{M}_1$  to  $\mathbf{M}_5$  mentioned above with a fixed common prior  $\theta_1$ .







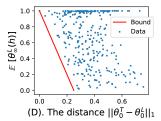


Figure 4. From left to right: (A). Rank 3,  $\mathbf{M}_3$  and  $\theta_1$ ,  $\theta_0^L$  is perturbed along six directions. (B). Rank 3,  $\mathbf{M}_3$  and  $\theta_1$ , sample  $\theta_0^L$  uniformly in  $\Delta^2$ . (C). Rank 4,  $\mathbf{M}_1'$  and  $\theta_1'$ , along 15 different directions. (D). Rank 4,  $\mathbf{M}_1'$  and  $\theta_1'$ , sample  $\theta_0^L$  uniformly in  $\Delta^3$ .

Let all matrices mentioned be column-normalized (this does not affect SCBI since cross-ratios and marginal conditions determines the Sinkhorn scaling results), we call the column determined by the true hypothesis h (the first column in our simulation) the target column ("tr. h" on Fig. 5), the column which  $\Re^s$  uses (argmin column) the relevant column ("rel. h") and the other column the irrelevant column ("irr. h"). Let  $\mathbf T$  be given, and let  $\mathbf L$  be obtained from  $\mathbf T$  by perturbing along the relevant / irrelevant column.

Without loss of generality, we assume that only one column of the learner's matrix L is perturbed at a time as other perturbations may be treated as compositions of such.

For each T and each column h', we apply 330 perturbations on concentric circles around T (the disc), 90 perturbations preserving the normalized-KL divergence ( $\mathrm{KL}(\mathbf{e}/n, \mathcal{N}_{\mathrm{vec}}(\mathbf{L}_{(\_,h')}/\mathbf{L}_{(\_,1)}, 1)$ ) used in  $\mathfrak{R}^{\mathrm{s}}$ ) from the target column and 50 linear interpolations with target column. Each point in Fig. 5 is estimated using a size- $10^4$  Monte Carlo method using Proposition 5.1. From the graphs, we can see that the successful rate varies continuously on perturbations, slow on one direction (the yellow strip crossing the center) and rapid on the perpendicular direction (color changed to blue rapidly).

## 6. Grid World: an Application

Consider a  $3 \times 5$  grid world with two possible terminal goals, A and B, and a starting position S as shown below. Let the reward at the terminal position  $h_t$  be R. Assuming no step costs, the value of a grid that distanced k from  $h_t$  is then  $R \times \gamma^k$  (in the RL-sense), where  $\gamma < 1$  is the discount factor.

A				В
		1		
	#	S	$\Rightarrow$	

Suppose the structure of the world is accessible to both agents whereas the true location of the goal  $h_t$  is only known to a teacher. The teacher performs a sequence of actions to teach  $h_t$  to a learner. At each round, there are three available actions, left, up and right. After observing the teacher's actions, the learner updates their belief on  $h_t$  accordingly.

We now compare BI and SCBI agents' behaviours under this grid world. In terms of previous notations, the hypothesis set  $\mathcal{H} = \{A, B\}$ , the data set  $\mathcal{D} = \{left, up, right\}$ . Let the learner's prior over  $\mathcal{H}$  be  $\theta_0 = (0.5, 0.5)$  and the true hypothesis  $h_t$  be A, then at each blue grid, agents'

(unnormalized) initial matrix 
$$\mathbf{M} = \begin{bmatrix} \text{left} & A & B \\ \text{left} & \gamma^{(k-1)} & \gamma^{(k+1)} \\ \text{up} & \gamma^{(k-1)} & \gamma^{(k-1)} \\ \gamma^{(k+1)} & \gamma^{(k-1)} \end{bmatrix}$$

Assume both BI teacher and SCBI teacher start with grid S. Based on M, the BI teacher would choose equally between *left* and up, whereas the SCBI teacher is more likely to choose *left* as the teacher's likelihood matrix  $\sum_{n=0}^{\infty} \binom{2/(3+3\gamma^2)}{2} \binom{2\gamma^2/(3+3\gamma^2)}{2} \binom{2\gamma^2}{3} \binom{2\gamma^2}{3$ 

 $\mathbf{T} = \begin{pmatrix} \frac{2/(3+3\gamma^2)}{1/3} & \frac{2\gamma^2/(3+3\gamma^2)}{1/3} \\ \frac{2\gamma^2/(3+3\gamma^2)}{2\gamma^2/(3+3\gamma^2)} & \frac{2}{2}(3+3\gamma^2) \end{pmatrix}, \text{ obtained from Sinkhorn scaling on M, assigns higher probability for } \textit{left.} \text{ Hence, comparing to the BI teacher who only aims for the final goal, the SCBI teacher tends to cooperate with the learner by selecting less ambiguous moves towards the goal. This point is aligned with the core idea of many existing models of cooperation in cognitive development (Jara-Ettinger et al., 2016; Bridgers et al., in press), pragmatic reasoning (Frank & Goodman, 2012; Goodman & Stuhlmüller, 2013) and robotics (Ho et al., 2016; Fisac et al., 2017).$ 

Moreover, even under the same teaching data, the SCBI learner is more likely to infer  $h_t$  than the BI learner. For instance, given the teacher's trajectory  $\{left, up\}$ , the left plot in Fig. 6 shows the SCBI and BI learners' posteriors on the true hypothesis  $h_t$ . Hence, comparing to the BI learner who reads the teacher's action literally, the SCBI learner interprets teacher's data corporately by updating belief sequentially after each round.

Regarding the stability, consider the case where the learner's discount factor is either greater or less (with equal probability) than the teacher's by 0.1. The right plot in Fig. 6 illustrates the expected difference between the learner's posterior on  $h_t$  after observing a teacher's trajectory of length 2 and the teacher's estimation of it.

As discussed in Sec 4.1, showing in Fig. 2, as the board gets wider and the number of possible goals gets more (i.e. the number of hypotheses increases), the gap between

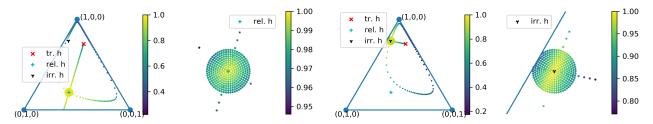


Figure 5. The perturbations on  $M_3$  along a column, and their zoomed-in version (with different color scale). The crosses shows the position of three normalized columns of  $T = M_3$ , the location of the dots represent the perturbed column of T (unperturbed columns are represented by crosses on figures which are not the center of disc) and whereas their colors depict the successful rate of inference. Left two figures are perturbations on the irrelevant column. Right two figures are perturbations on the relevant column.

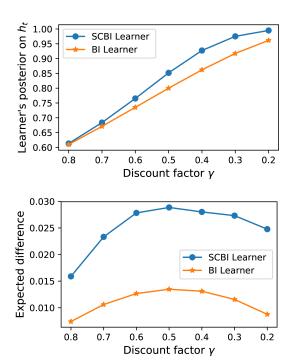


Figure 6. The top plot demonstrates that both BI and SCBI converge to the true hypothesis with SCBI having higher sample efficiency. The bottom plot shows that both BI and SCBI agents are robust to perturbations with SCBI relatively less stable.

posteriors of SCBI and BI learners will increase whereas the expected difference between agents for the same magnitude of perturbation will decrease. Thus, this example illustrates the consistency, sample efficiency, and stability of SCBI versus BI.

#### 7. Related Work

Literatures on Bayesian teaching (Eaves & Shafto, 2016; Eaves Jr et al., 2016), Rational Speech act theory (Frank & Goodman, 2012; Goodman & Stuhlmüller, 2013), and machine teaching (Zhu, 2015; 2013) consider the problem of selecting examples that improve a learner's chances of

inferring a concept. These literatures differ in that they consider the single step, rather than sequential problem, that they do not formalize learners who reason about the teacher's selection process, and that they models without a mathematical analysis.

The literature on pedagogical reasoning in human learning (Shafto & Goodman, 2008; Shafto et al., 2012; 2014) and cooperative inference (Yang et al., 2018; Wang et al., 2019a;b) in machine learning formalize full recursive reasoning from the perspectives of both the teacher and the learner. These only consider the problem of a single interaction between the teacher and learner.

The literature on curriculum learning considers sequential interactions with a learner by a teacher in which the teacher presents data in an ordered sequence (Bengio et al., 2009), and traces back to various literatures on human and animal learning (Skinner, 1958; Elman, 1993). Curriculum learning involves one of a number of methods for optimizing the sequence of data presented to the learner, most commonly starting with easier / simpler examples first and gradually moving toward more complex or less typical examples. Curriculum learning considers only problems where the teacher optimizes the sequence of examples, where the learner does not reason about the teaching.

#### 8. Conclusions

Cooperation is central to learning in humans and machines. We set out to provide a mathematical foundation for sequential cooperative Bayesian inference (SCBI). We presented new analytic results demonstrating the consistency and asymptotic rate of convergence of SCBI. Empirically, we demonstrated the sample efficiency and stability to perturbations as compared to Bayesian inference, and illustrated with a simple reinforcement learning problem. We therefore provide strong evidence that SCBI satisfies basic desiderata. Future work will aim to provide mathematical proofs of the empirically observed efficiency and stability.

### A. Proof of Consistency Theorems

## A.1. Proof of Theorem A.1

**Theorem A.1.** [Theorem 3.1,(Miescke & Liese, 2008, Theorem 7.115)] In BI, the sequence of posteriors  $(S_k)$  is strongly consistent at  $\widehat{\theta} = \delta_h$  for each  $h \in \mathcal{H}$ , with arbitrary choice of an interior point  $\theta_0 \in (\mathcal{P}(\mathcal{H}))^{\circ}$  (i.e.  $\theta_0(h) > 0$  for all  $h \in \mathcal{H}$ ) as prior.

*Proof.* We follow the same line as discussed right after this theorem in the paper. Let  $\theta_0 = (\theta_0(1), \theta_0(2), \dots, \theta_0(n))$  be the original prior, and let  $\theta_k = (\theta_k(1), \theta_k(2), \dots, \theta_k(m))$  be the posterior after having k data points  $d_1, d_2, \dots, d_k$ . Then for  $l \leq k$  and  $h \in \mathcal{H}$ , the posterior  $\theta_l(i) = \left(\mathcal{N}_{\text{ec}}(\operatorname{diag}(\mathbf{M}_{(d_l, \cdot)})\theta_{l-1})\right)(i)$  by Bayes' rule. In other words,

$$\theta_l(i) = \frac{\mathbf{M}_{(d_l,i)}[\theta_{(l-1)}(i)]}{\sum_{j=1}^m \mathbf{M}_{(d_l,j)}[\theta_{(l-1)}(j)]}.$$
 (5)

This is a recursive formula, so we may move forward to calculate  $\theta_l(i)$  from a smaller round index  $\theta_t(i)$  with t < l:

$$\theta_l(i) = \frac{\left[\prod_{s=t}^l \mathbf{M}_{(d_s,i)}\right] \theta_{(t-1)}(i)}{\sum_{j=1}^m \left[\prod_{s=t}^l \mathbf{M}_{(d_s,j)}\right] \theta_{(t-1)}(j)}.$$

This recursion stops at prior  $\theta_0$ , so we have an explicit expression of  $\theta_k$ :

$$\theta_k(i) = \frac{\left[\prod_{s=1}^k \mathbf{M}_{(d_s,i)}\right] \theta_0(i)}{\sum_{j=1}^m \left[\prod_{s=1}^k \mathbf{M}_{(d_s,j)}\right] \theta_0(j)}.$$
 (6)

It can be seen that for each hypothesis i, the denominator of the k-th posterior on i are the same, so we have

$$\frac{\theta_k(i)}{\theta_k(h)} = \frac{\left[\prod_{s=1}^k \mathbf{M}_{(d_s,h)}\right] \theta_0(i)}{\left[\prod_{s=1}^k \mathbf{M}_{(d_s,h)}\right] \theta_0(h)}.$$
 (7)

So we define  $\alpha_k(d)$  to be the frequency of the occurrence of data d in the first k rounds of a episode. And then

$$\log\left(\frac{\theta_k(i)}{\theta_k(h)}\right) = \log\left(\frac{\theta_0(i)}{\theta_0(h)}\right) + \sum_{d=1}^n \alpha_k(d) \log\left(\frac{\mathbf{M}_{(d,i)}}{\mathbf{M}_{(d,h)}}\right). \tag{8}$$

Since we know that the data  $(d_i)$  in the model is sampled following the i.i.d. with distribution  $\mathbf{M}_{(-,h)}$ , then for a fixed k,  $\alpha_k(i)$  follows the multinomial distribution with parameter  $\mathbf{M}_{(-,h)}$ .

By the strong law of large numbers,  $\frac{\alpha_k(i)}{k} \to \mathbf{M}_{(i,h)}$  almost surely as  $k \to \infty$ . Thus, when we rewrite the sample values

to random variable version,

$$\frac{1}{k} \log \left( \frac{\Theta_k(i)}{\Theta_k(h)} \right) \to \sum_{d=1}^n \mathbf{M}_{(d,h)} \log \left( \frac{\mathbf{M}_{(d,i)}}{\mathbf{M}_{(d,h)}} \right) \quad \text{a.s.}$$
(9)

That is,

$$\frac{1}{k}\log\left(\frac{\Theta_k(i)}{\Theta_k(h)}\right) \longrightarrow -\mathrm{KL}\left(\mathbf{M}_{(-,h)},\mathbf{M}_{(-,i)}\right) \quad \text{a.s.} \quad (10)$$

By the assumption in Section 2 of the paper that  $\mathbf{M}$  has distinct columns, the KL divergence between the i-th column and the h-th column is strictly positive, thus almost surely,  $\log\left(\frac{\Theta_k(i)}{\Theta_k(h)}\right) \to -\infty$ , or equivalently,  $\frac{\Theta_k(i)}{\Theta_k(h)} \to 0$ , for any  $i \neq h$ .

Therefore,  $\theta_k = (\theta_k(1), \theta_k(2), \dots, \theta_k(m)) \to \delta_h$  almost surely, equivalently, BI at  $\widehat{\theta}$  is strongly consistent.

#### A.2. Proof of Theorem 3.2

**Theorem A.2** (Theorem 3.2). In BI, with  $\widehat{\theta} = \delta_h$  for some  $h \in \mathcal{H}$ , let  $\Theta_k(h)(D_1, \ldots, D_k) := S_k(h|D_1, \ldots, D_k)$  be the h-component of posterior given  $D_1, \ldots, D_k$  as random variables valued in  $\mathcal{D}$ . Then  $\frac{1}{k}\log\left(\frac{\Theta_k(h)}{1-\Theta_k(h)}\right)$  converges to a constant  $\min_{h'\neq h}\left\{\mathrm{KL}(\mathbf{M}_{(.,h)},\mathbf{M}_{(.,h')})\right\}$  almost surely.

*Proof.* Follow the previous proof. First recall that  $\frac{1}{k}\log\left(\frac{\Theta_k(i)}{\Theta_k(h)}\right) \to -\mathrm{KL}(\mathbf{M}_{(\text{-},h)},\mathbf{M}_{(\text{-},i)}) \text{ almost surely.}$  Let  $\eta:=\mathrm{argmin}_{i\neq h}\left\{\mathrm{KL}(\mathbf{M}_{(\text{-},h)},\mathbf{M}_{(\text{-},i)})\right\}$ , then  $\Theta_k(\eta)$  decays slowest among  $\{\Theta_k(i):i\neq h\}$  almost surely.

Therefore, for the sample values  $\theta_k$ 's, asymptotically,

$$\frac{1}{k} \log \left[ \frac{\theta_k(\eta)}{\theta_k(h)} \right] \leq \frac{1}{k} \log \left[ \frac{1 - \theta_k(h)}{\theta_k(h)} \right] \leq \frac{1}{k} \log \left[ \frac{(m - 1)\theta_k(\eta)}{\theta_k(h)} \right].$$

So when we are taking limits  $k \to \infty$ , with probability one, we have

$$-\mathrm{KL}(\mathbf{M}_{(.,h)}, \mathbf{M}_{(.,\eta)}) \leq \lim_{k \to \infty} \frac{1}{k} \log \left[ \frac{1 - \theta_k(h)}{\theta_k(h)} \right]$$

$$\leq \lim_{k \to \infty} -\mathrm{KL}(\mathbf{M}_{(.,h)}, \mathbf{M}_{(.,\eta)}) + \frac{1}{k} \log(m-1)$$

$$= -\mathrm{KL}(\mathbf{M}_{(.,h)}, \mathbf{M}_{(.,\eta)}). \tag{11}$$

#### A.3. Proof of Theorem 3.5

To prove Theorem A.3, we need the following lemmas.

**Lemma A.3.1.** Given a fixed hypothesis  $h \in \mathcal{H}$ , for any  $\mu \in \mathcal{P}(\Delta^{m-1})$ ,

$$\mathbb{E}_{\mu}(\theta(h)) \le \mathbb{E}_{\Psi(h)(\mu)}(\theta(h)). \tag{12}$$

equality happens when  $\mathbf{M}^{\langle n\mathbf{x}\rangle}_{(i,h)} = \mathbf{M}^{\langle n\mathbf{x}\rangle}_{(j,h)}$  for any i,j and  $\mu$ -almost everywhere for  $\mathbf{x} \in \Delta^{m-1}$ .

Remark 6. This lemma shows that the expectation of  $\theta(h)$ , in each round is increasing, thus the sequence obtained from all the rounds has an limit since the sequence is monotonic and upper bounded by 1. To prove the theorem we, then just need to show the limit is 1.

*Proof.* We start from the right hand side of Eq. 12. Let  $\Delta$  denote  $\Delta^{m-1}$  for short.

$$\mathbb{E}_{\Psi(h)(\mu)}(\theta(h))$$

$$= \int_{\Delta} \theta(h) d(\Psi(h)(\mu))(\theta)$$

$$= \int_{\Delta} \sum_{d=1}^{n} \tau_{d}(T_{d}^{-1}(\theta)) \theta(h) d(T_{d*}(\mu))(\theta)$$

$$= \sum_{d=1}^{n} \int_{\Delta} \tau_{d}(\theta)(T_{d}(\theta))(h) d(T_{d*}(\mu))(T_{d}(\theta))$$

$$= \sum_{d=1}^{n} \int_{\Delta} \tau_{d}(\theta)(T_{d}(\theta))(h) d\mu(\theta)$$

$$= \sum_{d=1}^{n} \int_{\Delta} \frac{T_{d}(\theta)(h)}{n\theta(h)} T_{d}(\theta)(h) d\mu(\theta)$$

$$= \int_{\Delta} \sum_{d=1}^{n} \frac{T_{d}(\theta)(h)^{2}}{n\theta(h)} d\mu(\theta)$$

In the calculation, the bijectivity of  $T_d$  and the formula  $(T_{d*}(\mu))(E) = \mu(T_d^{-1}(E))$  is used (and will be used repetitively later).

Consider that by definition of the bijection  $T_d$ , the sum  $\sum_{d=1}^{n} T_d(\theta)(h) = n\theta(h)$  ( $T_d$  is the d-th row of Sinkhorn scaling by column sums  $n\theta$ ). Thus

$$\mathbb{E}_{\Psi(h)(\mu)}(\theta(h)) = \int_{\Delta} \frac{\sum_{d=1}^{n} T_{d}(\theta)(h)^{2}}{\sum_{d=1}^{n} T_{d}(\theta)(h)} d\mu(\theta)$$

$$\geq \int_{\Delta} \frac{\left(\sum_{d=1}^{n} T_{d}(\theta)(h)\right)^{2}}{n \sum_{d=1}^{n} T_{d}(\theta)(h)} d\mu(\theta)$$

$$= \int_{\Delta} \frac{1}{n} \sum_{d=1}^{n} T_{d}(\theta)(h) d\mu(\theta)$$

$$= \int_{\Delta} \theta(h) d\mu(\theta)$$

$$= \mathbb{E}_{\mu}(\theta(h)), \tag{13}$$

where  $\sum_{d=1}^{n} T_d(\theta)(h)^2 \geq \frac{1}{n} \left(\sum_{d=1}^{n} T_d(\theta)(h)\right)^2$  by Cauchy-Schwarz inequality, with equality achieved if and only if  $T_d(\theta)(h)$  is constant on d. Therefore, the equality of Eq. (13) is achieved when  $\mathbf{M}^{\langle n\mathbf{x}\rangle}_{(d,h)}$  is constant on d,  $\mu$ -almost everywhere for  $\mathbf{x} \in \Delta^{m-1}$ .

 $\theta(h) = 1 - \cdots$   $\theta(h) = b - \cdots$  Lemma A.3.4  $\theta(h) = a - \cdots$  Lemma A.3.2  $\theta(h) = 0 - \cdots$ 

Figure 7. Sketch of  $\Delta^{m-1}$ , for a general  $\theta$ , its y-coordinate is  $\theta(h)$ . The levels are compatible with proof of Theorem A.3. Lemma A.3.2 and Lemma A.3.4 are located where they contribute to prove the vanishing of measure in the limit.

The following lemmas helps showing that the measure  $\mu_k$  of the complement of a neighborhood of  $\delta_h \in \mathcal{P}(\mathcal{H})$  has limit 0.

**Lemma A.3.2.** Given M,  $h \in \mathcal{H}$  and prior  $\mu_0 \in \mathcal{P}(\Delta^{m-1})$  satisfying assumptions, we have

$$\mathbb{E}_{\mu_k}\left(\frac{\theta(h')}{\theta(h)}\right) = \mathbb{E}_{\mu_0}\left(\frac{\theta(h')}{\theta(h)}\right) \tag{14}$$

for any  $k \geq 0$  and any  $h' \neq h$ .

*Proof.* It suffices to prove the k=1 case for a general  $\mu_0$  (then we have the rest by induction).

$$\mathbb{E}_{\mu_{1}}\left(\frac{\theta(h')}{\theta(h)}\right) = \int_{\Delta} \left(\frac{\theta(h')}{\theta(h)}\right) d\mu_{1}(\theta)$$

$$= \int_{\Delta} \left(\frac{\theta(h')}{\theta(h)}\right) d(\Psi(h)(\mu_{0}))(\theta)$$

$$= \int_{\Delta} \sum_{d \in \mathcal{D}} \tau_{d}(T_{d}^{-1}(\theta)) \frac{\theta(h')}{\theta(h)} d(T_{d*}(\mu_{0}))(\theta)$$

$$= \sum_{d \in \mathcal{D}} \int_{\Delta} \tau_{d}(\theta) \frac{T_{d}(\theta)(h')}{T_{d}(\theta)(h)} d(T_{d*}(\mu_{0}))(T_{d}(\theta))$$

$$= \sum_{d \in \mathcal{D}} \int_{\Delta} \frac{T_{d}(\theta)(h)}{n\theta(h)} \frac{T_{d}(\theta)(h')}{T_{d}(\theta)(h)} d(\mu_{0})(\theta)$$

$$= \int_{\Delta} \sum_{d \in \mathcal{D}} \frac{T_{d}(\theta)(h')}{n\theta(h)} d(\mu_{0})(\theta)$$

$$= \int_{\Delta} \frac{\sum_{d \in \mathcal{D}} T_{d}(\theta)(h')}{n\theta(h)} d(\mu_{0})(\theta)$$

$$= \int_{\Delta} \frac{n\theta(h')}{n\theta(h)} d(\mu_{0})(\theta)$$

$$= \mathbb{E}_{\mu_{0}}\left(\frac{\theta(h')}{\theta(h)}\right). \tag{15}$$

**Lemma A.3.3.** The operator  $\Psi(h)$  preserves convex combinations of probability measures, i.e., for positive  $a_1, a_2, \ldots, a_l$  with  $\sum_{i=1}^l a_i = 1$  and probability measures  $\mu_1, \mu_2, \ldots, \mu_l$ ,

$$\Psi(h)\left(\sum_{i=1}^{l} a_i \mu_i\right) = \sum_{i=1}^{l} a_i \Psi(h)(\mu_i).$$

*Proof.* By definition, for any measurable set E in Borel  $\sigma$  algebra  $\mathfrak{A}$ ,

$$\Psi(h)(\mu)(E) := \int_{E} \sum_{d=1}^{n} \tau_{d}(T_{d}^{-1}(\theta)) d(T_{d*}(\mu))(\theta).$$

where every summand commutes with convex combination.

**Lemma A.3.4.** Given M, h, and  $\mu_0$  satisfying the assumptions, then for any 0 < a < b < 1,

$$\lim_{k \to \infty} \mu_k(\{\theta \in \Delta^{m-1} : a \le \theta(h) \le b\}) = 0 \tag{16}$$

*Proof.* We first show a property of  $\mu$  on the set  $\Delta_{[a,b]}:=\{\theta\in\Delta^{m-1}:a\leq\theta(h)\leq b\}.$ 

For any  $\mu$  supported on  $\Delta_{[a,b]}$  (that is,  $\mu(\Delta_{[a,b]})=1$ ), there is a positive number  $\epsilon_0$ , such that

$$\mathbb{E}_{\Psi^2(h)(\mu)}(\theta(h)) - \mathbb{E}_{\mu}(\theta(h)) \ge \epsilon_0. \tag{17}$$

According to the calculation in Lemma A.3.1, especially the first step of Eq. (13),

$$\mathbb{E}_{\Psi^{2}(h)(\mu)}(\theta(h))$$

$$= \int_{\Delta} \frac{\sum_{d=1}^{n} T_{d}(\theta)(h)^{2}}{\sum_{d=1}^{n} T_{d}(\theta)(h)} d(\Psi(h)(\mu))(\theta)$$

$$= \int_{\Delta} \sum_{e=1}^{n} \tau_{e}(T_{e}^{-1}(\theta)) \frac{\sum_{d=1}^{n} T_{d}(\theta)(h)^{2}}{\sum_{d=1}^{n} T_{d}(\theta)(h)} d(T_{e*}(\mu))(\theta)$$

$$= \int_{\Delta} \sum_{e=1}^{n} \tau_{e}(\theta) \frac{\sum_{d=1}^{n} T_{d}(T_{e}(\theta))(h)^{2}}{\sum_{d=1}^{n} T_{d}(T_{e}(\theta))(h)} d\mu(\theta)$$
(18)

Thus

$$\mathbb{E}_{\Psi^{2}(h)(\mu)}(\theta(h)) - \mathbb{E}_{\mu}(\theta(h))$$

$$= \int_{\Delta} \sum_{e=1}^{n} \tau_{e}(\theta) \frac{\sum_{d=1}^{n} T_{d}(T_{e}(\theta))(h)^{2}}{\sum_{d=1}^{n} T_{d}(T_{e}(\theta))(h)} - \theta(h) d\mu(\theta) \quad (19)$$

To show the claim, it suffices to find a positive lower bound of the integrand of Eq. (19),  $\Im(\theta):=\sum_{e=1}^n \tau_e(\theta) \frac{\sum_{d=1}^n T_d(T_e(\theta))(h)^2}{\sum_{d=1}^n T_d(T_e(\theta))(h)} - \theta(h)$ , for all  $\theta \in \Delta_{[a,b]}$ . Moreover, since  $\Delta_{[a,b]}$  is compact, we just need to show  $\Im(\theta) > 0$  on  $\Delta_{[a,b]}$ .

With Cauchy-Schwarz inequality used in Lemma A.3.1, we know

$$\mathfrak{I}(\theta) = \sum_{e=1}^{n} \tau_{e}(\theta) \frac{\sum_{d=1}^{n} T_{d}(T_{e}(\theta))(h)^{2}}{\sum_{d=1}^{n} T_{d}(T_{e}(\theta))(h)} - \theta(h)$$

$$\geq \sum_{e=1}^{n} \tau_{e}(\theta) \frac{1}{n} \left( \sum_{d=1}^{n} T_{d}(T_{e}(\theta))(h) \right) - \theta(h)$$

$$= \sum_{e=1}^{n} \tau_{e}(\theta) T_{e}(\theta)(h) - \theta(h)$$

$$= \sum_{e=1}^{n} \frac{T_{e}(\theta)(h)}{n\theta(h)} T_{e}(\theta)(h) - \theta(h)$$

$$\geq \frac{1}{n} T_{e}(\theta)(h) - \theta(h)$$

$$= \theta(h) - \theta(h) = 0$$
(20)

 $\Im(\theta)$  vanishes if and only if both line 2 and line 5 has equality, and we will discuss why these can not happen simultaneously.

The equality in line 5 requires that  $T_e(\theta)(h)$  are identical for all  $e \in \mathcal{D}$ , or more precisely, the vector  $\mathbf{M}^{\langle n\theta \rangle}_{(.,h)} = te_n$  has identical components. Further if equality in line 2 holds, the terms  $T_d(T_e(\theta))(h)$  are the same for all  $d \in \mathcal{D}$ . That is, by condition  $\sum_{e=1}^n T_e(\theta)(h) = n\theta(h)$  and  $\sum_{d=1}^n T_d(T_e(\theta))(h) = nT_e(\theta)(h)$ ,  $\mathfrak{I}(\theta)$  vanishes if and only if  $T_d(T_e(\theta))(h) = T_e(\theta)(h) = \theta(h)$  for all  $d, e \in \mathcal{D}$ .

We analyze the Sinkhorn scaled matrices in detail: Let  $\mathbf{M}^{\star} = \mathbf{M}^{\langle n\theta \rangle}$  be the scaled matrix whose e-th row is  $T_e(\theta)$ , and let  $\mathbf{M}^{(e)} = \mathbf{M}^{\langle nT_e(\theta) \rangle}$  be the scaled matrix whose d-th row is  $T_d(T_e(\theta))$ . Since  $\mathbf{M}^{\star}$  and each  $\mathbf{M}^{(e)}$  has the same h-th column, there are diagonal matrices  $\mathbf{D}^{(e)} = \mathrm{diag}(\frac{n\mathbf{M}_{(e,1)}^{\star}}{\sum_{i=1}^{n}\mathbf{M}_{(i,1)}^{\star}}, \frac{n\mathbf{M}_{(e,2)}^{\star}}{\sum_{i=1}^{n}\mathbf{M}_{(i,2)}^{\star}}, \dots, \frac{n\mathbf{M}_{(e,n)}^{\star}}{\sum_{i=1}^{n}\mathbf{M}_{(i,n)}^{\star}})$  such that  $\mathbf{M}^{(e)} = \mathbf{M}^{\star}\mathbf{D}^{(e)}$ . Since  $\mathbf{M}^{\star}$  and all  $\mathbf{M}^{(e)}$  are row-normalized to  $\mathbf{e}$  (i.e., their row sums are 1), we have the following equations from the row sums:

$$\mathfrak{S}(d,e) := \sum_{j=1}^{m} \mathbf{M}_{(d,j)}^{\star} \frac{n \mathbf{M}_{(e,j)}^{\star}}{\sum_{i=1}^{n} \mathbf{M}_{(i,j)}^{\star}} = 1$$
 (21)

for all  $d, e \in \mathcal{D}$  representing the d-th row-sum of  $\mathbf{M}^{(e)}$ .

Then we calculate  $(n-1)\sum_{e=1}^n\mathfrak{S}(e,e)-\sum_{d\neq e}\mathfrak{S}(d,e)$ . On the right hand side, since  $\mathfrak{S}(d,e)=1$  for every d,e, we have

$$(n-1)\sum_{e=1}^n \mathfrak{S}(e,e) - \sum_{d\neq e} \mathfrak{S}(d,e) = (n-1)n - (n^2 - n) = 0.$$

Meanwhile,

$$(n-1)\sum_{e=1}^{n}\mathfrak{S}(e,e) - \sum_{d\neq e}\mathfrak{S}(d,e)$$

$$= \sum_{j=1}^{m} \frac{n}{\sum_{i=1}^{n} \mathbf{M}_{(i,j)}^{\star}} \left(\sum_{e=1}^{n} (n-1)(\mathbf{M}_{(e,j)}^{\star})^{2} - \sum_{d\neq e} \mathbf{M}_{(d,j)}^{\star} \mathbf{M}_{(e,j)}^{\star}\right)$$

$$= \sum_{j=1}^{m} \frac{n}{\sum_{i=1}^{n} \mathbf{M}_{(i,j)}^{\star}} \left(\sum_{d < e} (\mathbf{M}_{(d,j)}^{\star} - \mathbf{M}_{(e,j)}^{\star})^{2}\right)$$

$$= 0$$

$$(22)$$

Therefore,  $\mathbf{M}_{(d,j)}^{\star} = \mathbf{M}_{(e,j)}^{\star}$  for any d,e, and j. Therefore, the rows of  $\mathbf{M}^{\star}$  are identical, so the columns of  $\mathbf{M}^{\star}$  are all parallel (or say, collinear as vectors, i.e. one is a scalar-multiple of the other) to each other.

By Sinkhorn scaling theory (Fienberg et al., 1970), the crossratios are invariant. Since M is a positive matrix and has distinct (non-parallel) columns, the  $2 \times 2$  cross-ratios are not identically 1, however,  $\mathbf{M}^{\star}$  — a scaled matrix of  $\mathbf{M}$  — has cross ratios identically 1. Therefore our assumption that  $\mathfrak{I}(\theta)=0$  cannot happen, and by compactness of  $\Delta_{[a,b]}$  and continuity of  $\mathfrak{I}(\theta)$ , we can conclude that  $\mathfrak{I}(\theta)$  has a lower bound  $\epsilon_0>0$  on  $\Delta_{[a,b]}$ .

Therefore,

$$\mathbb{E}_{\Psi^{2}(h)(\mu)}(\theta(h)) - \mathbb{E}_{\mu}(\theta(h))$$

$$= \int_{\Delta} \Im(\theta) d\mu(\theta)$$

$$\geq \int_{\Delta} \epsilon_{0} d\mu(\theta)$$

$$= \epsilon_{0} \mu(\Delta) = \epsilon_{0}. \tag{23}$$

Thus we prove the property Eq. (17).

We prove the lemma by contradiction:

Suppose the limit does not exist or the limit is nonzero. In either case, there exists a positive real number  $\epsilon > 0$ , such that there are infinitely many integers, or say a sequence  $(k_i)$  such that

$$\mu_{k_i}(\{a \le \theta(h) \le b\}) > \epsilon.$$

We may assume  $k_i$  contains no consecutive elements, i.e.,  $k_{i+1}-k_i>1$  for all i, otherwise, we can always find a subsequence satisfying this (for example, choose the sequence of all odd or even  $k_i$ 's, at least one of them is infinite, so we have a sequence).

For a  $\mu$ -measurable set E, let  $\mu|_E$  be the restriction of  $\mu$  on E, which can be treated as a measure on  $\Delta$  by setting the measure of the complement  $E^c$  zero (but the measure of  $\Delta$  is no longer 1). We scale it to  $\widehat{\mu}|_E:=(\mu(E))^{-1}\mu|_E$  to make it a probability measure, then  $\mu_{k_i}=[\mu_{k_i}(\Delta_{[a,b]})]\widehat{\mu}_{k_i}|_{\Delta_{[a,b]}}+[1-\mu_{k_i}(\Delta_{[a,b]})]\widehat{\mu}_{k_i}|_{(\Delta-\Delta_{[a,b]})}.$  Thus according to Lemma A.3.3,

$$\mathbb{E}_{\Psi^{2}(h)(\mu_{k_{i}})}(\theta(h))$$

$$= \mu_{k_{i}}(\Delta_{[a,b]})\mathbb{E}_{\Psi^{2}(h)(\widehat{\mu}_{k_{i}}|_{\Delta_{[a,b]}})}(\theta(h))$$

$$+(1-\mu_{k_{i}}(\Delta_{[a,b]}))\mathbb{E}_{\Psi^{2}(h)(\widehat{\mu}_{k_{i}}|_{(\Delta-\Delta_{[a,b]})})}(\theta(h))$$

$$\geq \mu_{k_{i}}(\Delta_{[a,b]})(\mathbb{E}_{\widehat{\mu}_{k_{i}}|_{\Delta_{[a,b]}}}(\theta(h))+\epsilon_{0})$$

$$+(1-\mu_{k_{i}}(\Delta_{[a,b]}))\mathbb{E}_{\widehat{\mu}_{k_{i}}|_{(\Delta-\Delta_{[a,b]})}}(\theta(h))$$

$$\geq \epsilon\epsilon_{0}+\mathbb{E}_{\mu_{k_{i}}}(\theta(h))$$
(24)

By Lemma A.3.1, we can see that  $\mathbb{E}_{\mu_{k+1}}[\theta(h)] \geq \mathbb{E}_{\mu_k}[\theta(h)]$ , and there is a sequence  $(k_i)$  such that  $\mathbb{E}_{\mu_{k_i+2}}[\theta(h)] \geq \mathbb{E}_{\mu_k}[\theta(h)] + \epsilon_0 \epsilon$ . Thus  $\mathbb{E}_{\mu_{k_i+2}}[\theta(h)] \geq \mathbb{E}_{\mu_0}[\theta(h)] + i\epsilon_0 \epsilon$ , so  $\lim_{k \to \infty} \mathbb{E}_{\mu_k}[\theta(h)] = \infty$ .

However,  $\theta(h) \leq 1$ , we have  $\mathbb{E}_{\mu_k}(\theta(h)) \leq 1$  for all k, which is a contradiction. Therefore, we know that Eq. (16) holds.

**Theorem A.3** (Theorem 3.5). In SCBI, let M be a positive matrix. If the teacher is teaching one hypothesis h (i.e.,  $\widehat{\theta} = \delta_h \in \mathcal{P}(\mathcal{H})$ ), and the prior distribution  $\mu_0 \in \mathcal{P}(\Delta^{m-1})$  satisfies  $\mu_0 = \delta_{\theta_0}$  with  $\theta_0(h) > 0$ , then the estimator sequence  $(S_k)$  is consistent, for each  $h \in \mathcal{H}$ , i.e., the posterior random variables  $(\Theta_k)_{k \in \mathbb{N}}$  converge to the constant random variable  $\widehat{\theta}$  in probability.

Some notions used in the proof are visualized in Fig. 7.

*Proof.* Let  $Z_0$  be a random variable with sample space  $\Delta^{m-1}$  such that the  $\text{Law}(Z_0) = \mu_0$ . This is the initial state in SCBI. The posteriors in the following rounds are determined by the sequence of data taught by teacher, which makes the posteriors random variables as well. Let  $Z_k$  be the random variable representing the posterior after k-rounds of SCBI, the law of  $Z_k$  is given by  $\text{Law}(Z_k) = \mu_k = |\Psi(h)|^k (\mu_0)$  according to the definition of  $\Psi(h)$ .

The consistency mentioned in the theorem is equivalent to that the sequence  $(Z_k)$  converges to  $\widehat{Z}$  with  $\text{Law}(\widehat{Z}) = \widehat{\mu}$  in probability where  $\widehat{\mu} = \delta_{\widehat{\theta}}$ .

We prove the theorem by contradiction. Suppose  $Z_k \to \widehat{Z}$  in probability is not valid, i.e., there exists  $\epsilon > 0$  such that

$$\lim_{k \to \infty} \Pr(d(Z_k, \widehat{Z}) > \epsilon)$$
 (25)

does not exist or the limit is positive, where the metric d on  $\Delta^{m-1}$  is the Euclidean distance inherited from  $\mathbb{R}^m$ . In either case, there is a real number C>0 such that

$$\Pr(d(Z_{k'}, \widehat{Z}) > \epsilon) > C \tag{26}$$

for a subsequence  $(Z_{k'})$  of  $(Z_k)$ .

Let  $R:=\mathbb{E}_{\mu_0}\left[\frac{1-\theta(h)}{\theta(h)}\right]=\frac{1-\theta_0(h)}{\theta_0(h)}$ , let  $a=\frac{1}{4R/C+1}$  and  $b=1-\epsilon$ . By Lemma A.3.4, there exists N>0 such that for all k>N,

$$\mu_k(\Delta_{[a,b]}) < C/2.$$

Therefore, for all the terms in (k') satisfying k' > N,  $\mu_{k'}(\{\theta:\theta(h)< a\}) > C/2$ . Furthermore,

$$\mathbb{E}_{\mu_{k'}} \left[ \frac{1 - \theta(h)}{\theta(h)} \right]$$

$$\geq \int_{\{\theta:\theta(h) < a\}} \left[ \frac{1 - \theta(h)}{\theta(h)} \right] d\mu_{k'}(\theta)$$

$$\geq \left[ \frac{1 - a}{a} \right] \frac{C}{2}$$

$$= \left[ \frac{4R}{C} \right] \frac{C}{2}$$

$$= 2R > R. \tag{27}$$

However, by Lemma A.3.2,

$$\mathbb{E}_{\mu_{k'}} \left[ \frac{1 - \theta(h)}{\theta(h)} \right]$$

$$= \mathbb{E}_{\mu_{k'}} \left[ \frac{\sum_{h' \neq h} \theta(h')}{\theta(h)} \right]$$

$$= \sum_{h' \neq h} \mathbb{E}_{\mu_{k'}} \left[ \frac{\theta(h')}{\theta(h)} \right]$$

$$= \sum_{h' \neq h} \mathbb{E}_{\mu_0} \left[ \frac{\theta(h')}{\theta(h)} \right]$$

$$= \mathbb{E}_{\mu_0} \left[ \frac{\sum_{h' \neq h} \theta(h')}{\theta(h)} \right]$$

$$= \mathbb{E}_{\mu_0} \left[ \frac{1 - \theta(h)}{\theta(h)} \right]$$

$$= R, \tag{28}$$

which is a contradiction. Therefore,

$$\lim_{k \to \infty} \Pr(d(Z_k, \widehat{Z}) > \epsilon) = 0.$$
 (29)

And the sequence of SCBI estimators is consistent at  $\widehat{\theta}$ .

## A.4. Proof of Theorem 3.6

**Theorem A.4** (Theorem 3.6). With matrix  $\mathbf{M}$ , hypothesis  $h \in \mathcal{H}$ , and a prior  $\mu_0 = \delta_{\theta_0} \in \mathcal{P}(\Delta^{m-1})$  same as in Theorem. A.3, let  $\theta_k$  denote a sample value of the posterior  $\Theta_k$  after k rounds of SCBI, then

$$\lim_{k \to \infty} \mathbb{E}_{\mu_k} \left[ \frac{1}{k} \log \left( \frac{\theta_k(h)}{1 - \theta_k(h)} \right) \right] = \mathfrak{R}^s(\mathbf{M}; h)$$
 (30)

where  $\mathfrak{R}^{s}(\mathbf{M};h) := \min_{h \neq h'} \mathrm{KL}\left(\mathbf{M}_{(.,h)}^{\sharp}, \mathbf{M}_{(.,h')}^{\sharp}\right)$  with  $\mathbf{M}^{\sharp} = \mathscr{N}_{col}(\mathrm{diag}(\mathbf{M}_{(.,h)})^{-1}\mathbf{M})$ . Thus we call  $\mathfrak{R}^{s}(\mathbf{M};h)$  the asymptotic rate of convergence (RoC) of SCBI.

*Proof.* We treat  $\theta_k$  as random variables, then

$$\mathbb{E}_{\mu_{k+1}}\!\!\left(\log\left[\frac{\theta_{k+1}(h')}{\theta_{k+1}(h)}\right]\right) = \mathbb{E}_{\mu_{k}}\!\left(\log\left[\frac{\theta_{k}(h')}{\theta_{k}(h)}\right]\right) + W_{k}^{h'},$$

where

$$W_k^{h'} = -\mathbb{E}_{\mu_k}\left[\mathrm{KL}\left(\mathscr{N}_{\mathrm{vec}}(\mathbf{M}^{\langle n\theta_k\rangle}_{(.,h)}), \mathscr{N}_{\mathrm{vec}}(\mathbf{M}^{\langle n\theta_k\rangle}_{(.,h')})\right)\right].$$

We can get it from the following calculation ( $\Delta$  represents

the simplex  $\Delta^{m-1}$ ):

$$\mathbb{E}_{\mu_{k+1}} \left( \log \left[ \frac{\theta_{k+1}(h')}{\theta_{k+1}(h)} \right] \right) \\
= \int_{\Delta} \log \left[ \frac{\theta(h')}{\theta(h)} \right] d\mu_{k+1}(\theta) \\
= \int_{\Delta} \sum_{d=1}^{n} \tau_{d}(T_{d}^{-1}(\theta)) \log \left[ \frac{\theta(h')}{\theta(h)} \right] d(T_{d*}(\mu_{k+1}))(\theta) \\
= \sum_{d=1}^{n} \int_{\Delta} \tau_{d}(\theta) \log \left[ \frac{T_{d}(\theta)(h')}{T_{d}(\theta)(h)} \right] d(\mu_{k})(\theta) \\
= \int_{\Delta} \sum_{d=1}^{n} \frac{T_{d}(\theta)(h)}{n\theta(h)} \log \left[ \frac{T_{d}(\theta)(h')}{T_{d}(\theta)(h)} \right] d(\mu_{k})(\theta) \\
= \int_{\Delta} \sum_{d=1}^{n} \frac{T_{d}(\theta)(h)}{n\theta(h)} \left\{ \log \left[ \frac{T_{d}(\theta)(h')}{n\theta(h')} \frac{n\theta(h)}{T_{d}(\theta)(h)} \right] + \log \left[ \frac{n\theta(h')}{n\theta(h)} \right] \right\} d(\mu_{k})(\theta) \\
= \int_{\Delta} -KL \left( \mathcal{N}(\mathbf{M}_{(\cdot,h)}^{(n\theta)}), \mathcal{N}(\mathbf{M}_{(\cdot,h')}^{(n\theta)}) \right) d\mu_{k} \\
+ \int_{\Delta} \log \left[ \frac{\theta(h')}{\theta(h)} \right] d\mu_{k} \\
= W_{k}^{h'} + \mathbb{E}_{\mu_{k}} \left( \log \left[ \frac{\theta_{k}(h')}{\theta_{k}(h)} \right] \right). \tag{31}$$

Next, we show

$$\lim_{k \to \infty} \mathbb{E}_{\mu_k} \frac{1}{k} \log \left( \frac{\theta_k(h')}{\theta_k(h)} \right) = -\text{KL}\left( \mathbf{M}_{(-,h)}^{\sharp}, \mathbf{M}_{(-,h')}^{\sharp} \right), \quad (32)$$

and then by a similar argument in the proof of Theorem A.2, we can show the result in this theorem.

To show Eq. (32), we can make use of Eq. (31). By showing that  $W_k^{h'}$  converges to  $-\mathrm{KL}\left(\mathbf{M}_{(-,h)}^\sharp,\mathbf{M}_{(-,h')}^\sharp\right)$ , we can conclude that  $\mathbb{E}_{\mu_k}\frac{1}{k}\log\left(\frac{\theta_k(h')}{\theta_k(h)}\right)$ , as the average of  $(W_i^{h'})$  on the first k-terms, converges to  $-\mathrm{KL}\left(\mathbf{M}_{(-,h)}^\sharp,\mathbf{M}_{(-,h')}^\sharp\right)$  as well.

To prove this, we need the following result from direct calculation:

**Lemma A.4.1.** Given a  $n \times 2$  positive matrix  $[\mathbf{a}, \mathbf{b}]$  with columns as n-vectors  $\mathbf{a} = (a_1, a_2, \dots, a_n)^{\top}$  and  $\mathbf{b} = (b_1, b_2, \dots, b_n)^{\top}$  with  $\sum_{i=1}^n a_i = \sum_{i=1}^n b_i = 1$ , consider the  $2 \times 2$  cross-ratios:  $C_i := CR(1, 2; 1, i) = \frac{a_1b_i}{a_ib_1}$ , then  $\mathrm{KL}(\mathbf{a}, \mathbf{b}) = \log \left(\sum_{i=1}^n a_i C_i\right) - \sum_{i=1}^n a_i \log C_i$ . With fixed  $C_i \in (0, \infty)$  for  $i = 1, 2, \dots, n$ ,  $\mathrm{KL}(\mathbf{a}, \mathbf{b})$  is continuous and bounded about  $\mathbf{a} \in \Delta^{n-1}$ .

*Proof of Lemma A.4.1.* The formula of  $KL(\mathbf{a}, \mathbf{b})$  is from direct calculation.

The KL-divergence is continuous and bounded since by the formula, every part is continuous and bounded given the restrictions on a and  $C_i$ .

Now we continue to prove Theorem A.4:

By continuity of the KL-divergence given fixed cross-ratios, for any  $\epsilon > 0$ , we find a number  $\delta > 0$  such that for any  $\theta \in \Delta^{m-1}$  with  $\theta(h) > 1 - \delta$ ,

$$\left| \mathrm{KL} \left( \mathcal{N}_{\mathrm{Vec}}(\mathbf{M}_{(-,h)}^{\langle n\theta \rangle}), \mathcal{N}_{\mathrm{Vec}}(\mathbf{M}_{(-,h')}^{\langle n\theta \rangle}) \right) - \mathrm{KL} \left( \mathbf{M}_{(-,h)}^{\sharp}, \mathbf{M}_{(-,h')}^{\sharp} \right) \right| < \frac{\epsilon}{3}. \tag{33}$$

Further, according to Theorem A.3, and the boundedness from Lemma A.4.1, we can find a number N>0, such that for any k>N, we have  $\mu_k(\{\theta:\theta(h)<1-\delta\})< C$  where C satisfies

$$C \cdot \sup_{\theta \in \Delta^{m-1}} \left\{ \mathrm{KL} \left( \mathscr{N}_{\mathrm{vec}}(\mathbf{M}_{(.,h)}^{\langle n\theta \rangle}), \mathscr{N}_{\mathrm{vec}}(\mathbf{M}_{(.,h')}^{\langle n\theta \rangle}) \right) \right\} < \frac{\epsilon}{3}. \quad (34)$$

The expectation  $W_k^{h'}$  can be split into two parts,  $W_k^{h'} = -W_> - W_<$  where

$$W_{>} = \int_{\theta(h)>1-\delta} KL\left(\mathcal{N}_{\text{vec}}(\mathbf{M}_{(-,h)}^{(n\theta)}), \mathcal{N}_{\text{vec}}(\mathbf{M}_{(-,h')}^{(n\theta)})\right) d\mu_k(\theta), \quad (35)$$

and

$$W_{\leq} = \int_{\theta(h) \leq 1-\delta} \mathrm{KL}\left(\mathcal{N}_{\text{ec}}(\mathbf{M}_{(-,h)}^{\langle n\theta \rangle}), \mathcal{N}_{\text{vec}}(\mathbf{M}_{(-,h')}^{\langle n\theta \rangle})\right) d\mu_k(\theta). \tag{36}$$

Similarly, since  $\mu_k$  is a probability measure  $\mathrm{KL}\left(\mathbf{M}_{(.,h)}^\sharp,\mathbf{M}_{(.,h')}^\sharp\right)=K_>+K_<$  where

$$K_{>} = \int_{\theta(h)>1-\delta} \mathrm{KL}\left(\mathbf{M}_{(-,h)}^{\sharp}, \mathbf{M}_{(-,h')}^{\sharp}\right) \mathrm{d}\mu_{k}(\theta), \tag{37}$$

and

$$K_{\leq} = \int_{\theta(h) \leq 1-\delta} \mathrm{KL}\left(\mathbf{M}_{(.,h)}^{\sharp}, \mathbf{M}_{(.,h')}^{\sharp}\right) \mathrm{d}\mu_k(\theta). \tag{38}$$

Then we have

$$\left| W_k^{h'} + \text{KL}\left(\mathbf{M}_{(-,h')}^{\sharp}, \mathbf{M}_{(-,h')}^{\sharp}\right) \right| \le |K_{>} - W_{>}| + |W_{<}| + |K_{<}|.$$
 (39)

The choice of  $\delta$  can make a good estimate of the integral on  $\theta(h) > 1 - \delta$ .

$$|K_{>} - W_{>}|$$

$$\leq \frac{\epsilon}{3}(1 - C)$$

$$< \frac{\epsilon}{3}.$$
(40)

For the other two terms, directly from condition Eq. (34), we have  $|K_<|<\frac{\epsilon}{3}$  and  $|W_<|<\frac{\epsilon}{3}$ , and hence  $|K_>-W_>|+|K_<|+|W_<|<\epsilon$ .

Therefore,  $W_k^{h'}$  converges to  $-\mathrm{KL}\left(\mathbf{M}_{(.,h)}^\sharp,\mathbf{M}_{(.,h')}^\sharp\right)$ .

#### A.5. An Example on a 2 by 2 Matrix

Let  $\mathcal{H}=\{h_1,h_2\}$ ,  $\mathcal{D}=\{d_1,d_2\}$ , and the shared joint distribution be  $\mathbf{M}^{JD}=\begin{pmatrix} d_1 & 0.3 & 0.3 \\ d_2 & \begin{pmatrix} 0.3 & 0.3 \\ 0.1 & 0.3 \end{pmatrix}$ . Further assume

that the learner has uniform prior on  $\mathcal{H}$ , i.e.  $S_0 = \theta_0 =$ (0.5, 0.5) and the true hypothesis given to the teachers is  $h_1$ . In round 1, the BI teacher will sample a data from  $\mathcal{D}$ according to the first column of  $\mathbf{M} = \begin{pmatrix} 0.75 & 0.5 \\ 0.25 & 0.5 \end{pmatrix}$ , which is obtained by column normalizing  $\mathbf{M}^{JD}$ . On the contrast, the SCBI teacher will form his likelihood matrix by first doing  $(\mathbf{r}_1, \mathbf{c}_1)$ -Sinkhorn scaling on M, then column normalization if needed, where  $\mathbf{r}_1=(1,1)$  and  $\mathbf{c}_1=(1,1)$  based on the uniform priors. The resulting limit matrix (with precision of three decimals) is  $\mathbf{M}_1^* = \begin{pmatrix} 0.634 & 0.366 \\ 0.366 & 0.634 \end{pmatrix}$ , which is already column normalized. Hence the SCBI teacher will teach according to the first column of the  $M_1 = M_1^*$ . Suppose that  $d_1$  is sampled by both teachers. The posterior for the BI learner is  $S_1^{b}(d_1) = (0.6, 0.4)$  (normalizing the first row of M). The posterior for the SCBI learner is  $S_1^{\rm s}(d_1) =$ (0.643, 0.366) (the first row of  $M_1$ ).

Similarly, in round 2, the SCBI teacher would update his likelihood matrix by first doing  $(\mathbf{r}_2,\mathbf{c}_2)\text{-Sinkhorn scaling}$  on  $\mathbf{M}_1,$  where  $\mathbf{r}_2=(1,1)$  and  $\mathbf{c}_2=(0.643,0.366)\times 2=(1.268,0.732).$  The resulting limit matrix is  $\mathbf{M}_2^*=\begin{pmatrix} 0.758 & 0.242 \\ 0.51 & 0.49 \end{pmatrix}$ . Then through column normalizing  $\mathbf{M}_2^*,$  a

updated likelihood matrix  $\mathbf{M}_2 = \begin{pmatrix} 0.60 & 0.33 \\ 0.4 & 0.67 \end{pmatrix}$  is obtained. The SCBI teacher will teach according to the first column of the  $\mathbf{M}_2$ . Whereas the BI teacher will again sample another data according to the the first column of  $\mathbf{M}$ . Suppose that  $d_1$  is sampled for both teachers. The posteriors for BI and SCBI learners are  $S_2^{\mathrm{b}}(d_1,d_1)=(0.692,0.308)$  and  $S_1^{\mathrm{s}}(d_1,d_1)=(0.758,0.242)$  respectively.

Although same teaching points are assumed, the SCBI learner's posterior on the true hypothesis  $h_1$  is higher than the BI learner in both rounds. Moreover, notice that the KL divergence between  $h_1$  and  $h_2$  is increasing as the likelihood matrix is updating through the SCBI. This will eventually lead much faster convergence for the SCBI learner.

#### **B.** Calculations about Sample Efficiency

Here we compute the expectation  $\mathfrak E$  mentioned in Sec. 4.1 of the paper, for matrices of size  $n \times 2$ .

We first calculate the average of RoC for a particular matrix

M. For simplicity, let  $\mathbf{x} = \mathbf{M}_{(\underline{\ },1)}$  and  $\mathbf{y} = \mathbf{M}_{(\underline{\ },2)}$ .

$$\frac{1}{2} \sum_{h=1}^{2} \mathfrak{R}^{b}(\mathbf{M}; h)$$

$$= \frac{1}{2} \left( \text{KL}(\mathbf{M}_{(.,1)}, \mathbf{M}_{(.,2)}) + \text{KL}(\mathbf{M}_{(.,2)}, \mathbf{M}_{(.,1)}) \right)$$

$$= \frac{1}{2} \left( \text{KL}(\mathbf{x}, \mathbf{y}) + \text{KL}(\mathbf{y}, \mathbf{x}) \right)$$

$$= \frac{1}{2} \left( \sum_{i=1}^{n} (\mathbf{x}_{i} - \mathbf{y}_{i}) (\ln \mathbf{x}_{i} - \ln \mathbf{y}_{i}) \right). \tag{41}$$

To calculate that for SCBI, denote  $\mathbf{x}/\mathbf{y} = \mathscr{N}_{vec}(\mathbf{v})$  the normalization of vector  $\mathbf{v}$  with  $\mathbf{v}_i = \mathbf{x}_i/\mathbf{y}_i$ .

$$\frac{1}{2} \sum_{h=1}^{2} \mathfrak{R}^{s}(\mathbf{M}; h)$$

$$= \frac{1}{2} \left[ \operatorname{KL} \left( \frac{\mathbf{e}}{n}, \mathbf{x}/\mathbf{y} \right) + \operatorname{KL} \left( \frac{\mathbf{e}}{n}, \mathbf{y}/\mathbf{x} \right) \right]$$

$$= \frac{1}{2} \left[ \frac{1}{n} \sum_{i=1}^{n} \left( -2 \ln n - \ln \frac{\mathbf{x}_{i}}{\mathbf{y}_{i}} + \ln \left( \sum_{j=1}^{n} \frac{\mathbf{x}_{j}}{\mathbf{y}_{j}} \right) \right) - \ln \frac{\mathbf{y}_{i}}{\mathbf{x}_{i}} + \ln \left( \sum_{j=1}^{n} \frac{\mathbf{y}_{j}}{\mathbf{x}_{j}} \right) \right) \right]$$

$$= \frac{1}{2} \left[ \ln \left( \sum_{j=1}^{n} \frac{\mathbf{x}_{j}}{\mathbf{y}_{j}} \right) + \ln \left( \sum_{j=1}^{n} \frac{\mathbf{y}_{j}}{\mathbf{x}_{j}} \right) \right] - \ln n. \quad (42)$$

The simulation of  $\mathfrak{P}$  is based on the above calculations. For  $\mathfrak{E}$ , the above expressions can be further simplified.

Given  $\mathbf{M} = (\mathbf{x}, \mathbf{y})$  uniformly distributed in  $(\Delta^{n-1})^2$ , with measure  $\nu \otimes \nu$  where  $\nu$  is the measure of uniform probability distribution on  $\Delta^{n-1}$ , we can calculate the expected value,

$$\mathbb{E}\left[\frac{1}{2}\sum_{h=1}^{2}\mathfrak{R}^{b}(\mathbf{M};h)\right]$$

$$=\frac{1}{2}\int_{(\Delta^{n-1})^{2}}\sum_{i=1}^{n}(\mathbf{x}_{i}-\mathbf{y}_{i})(\ln\mathbf{x}_{i}-\ln\mathbf{y}_{i})d\nu(\mathbf{x})d\nu(\mathbf{y})$$

$$=n\int_{\Delta^{n-1}}\mathbf{x}_{1}\ln\mathbf{x}_{1}d\nu(\mathbf{x})-n\int_{(\Delta^{n-1})^{2}}\mathbf{x}_{1}\ln\mathbf{y}_{1}d\nu(\mathbf{x})d\nu(\mathbf{y})$$

$$=n\int_{0}^{1}x(n-1)(1-x)^{n-2}\ln xdx+$$

$$n\int_{0}^{1}\int_{0}^{1}x(n-1)^{2}(1-x)^{n-2}(1-y)^{n-2}\ln ydxdy$$

$$=\frac{n-1}{n}.$$
(43)

Here we use the fact that

$$\int_{\{\theta \in \Delta^{n-1}: \theta(h) = a\}} d\mathbf{x} = (n-1)(1-a)^{n-2}.$$

Furthermore, since integral on  $(\Delta^{n-1})^2$  with measure  $\nu \otimes \nu$  is symmetric on  $\mathbf{x}$  and  $\mathbf{y}$ , we have

$$\mathfrak{E} = \int_{(\Delta^{n-1})^2} \ln \left( \sum_{i=1}^n \frac{\mathbf{x}_i}{\mathbf{y}_i} \right) d\nu(\mathbf{x}) d\nu(\mathbf{y}) - \ln n - \frac{n-1}{n}.$$
 (44)

In general, we calculate the integral in Eq. (44) by Monte Carlo method since other numerical integral methods we tried becomes slow dramatically as n grows. In particular, when n=2, an expression related to the dilogarithm  $Li_2$  can be obtained (can be easily checked in Wolfram software).

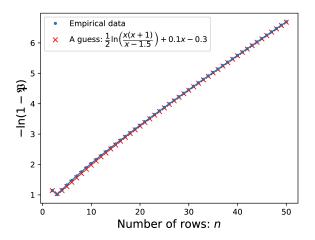


Figure 8. A guess of the  $\mathfrak P$  values for each n together with empirical data

And we have an empirical formula to describe the relation between  $\mathfrak{P}$  and n, shown in Fig. 8:

$$\mathfrak{P}_{(n,2)} = 1 - \sqrt{\frac{x - 1.5}{x(x+1)}} e^{-0.1x + 0.3}$$
 (45)

## C. Theory and Empirical Data on Stability

#### C.1. Proof of Proposition 5.1

**Proposition C.1** (Proposition 5.1). Given a sequence of identical independent  $\mathcal{D}$ -valued random variables  $(D_i)_{i\geq 1}$  following uniform distribution. Let  $\mu_0 \in \mathcal{P}(\Delta^{m-1})$  be a prior distribution on  $\Delta^{m-1}$ , and  $\mu_{k+1} = \Psi^{\mathbf{L}}_{D_k}(\mu_k)$ , then  $\mu_k$  converges, in probability, to  $\sum_{i\in\mathcal{H}} a_i\delta_i$ , where  $a_i = \mathbb{E}_{\mu_0} [\theta(i)]$ .

To show the above proposition, we need the following lemma:

**Lemma C.1.1.** Given the conditions in Proposition C.1, then for any  $k \in \mathbb{N}$  and  $h \in \mathcal{H}$ ,

$$\mathbb{E}_{\mu_k}(\theta(h)) = \mathbb{E}_{\mu_0}(\theta(h)). \tag{46}$$

*Proof.* It suffices to show  $\mathbb{E}_{\mu_{k+1}}(\theta(h)) = \mathbb{E}_{\mu_k}(\theta(h))$  for any k.

$$\mathbb{E}_{\mu_{k+1}}(\theta(h)) = \int_{\Delta^{m-1}} \theta(h) d(\mu_{k+1})(\theta)$$

$$= \sum_{d \in \mathcal{D}} D_k(d) \int_{\Delta^{m-1}} T_d(\theta)(h) d\mu_k(\theta)$$

$$= \int_{\Delta^{m-1}} \frac{1}{n} \sum_{d \in \mathcal{D}} T_d(\theta)(h) d(\mu_k)(\theta)$$

$$= \int_{\Delta^{m-1}} \theta(h) d(\mu_k)(\theta) = \mathbb{E}_{\mu_k}(\theta(h)).$$

*Proof of Proposition C.1.* We first show the following result:

For any  $\epsilon > 0$ , let

$$\Delta_{\epsilon} := \left\{ \theta \in \Delta^{m-1} : \theta(i) \le 1 - \epsilon, \forall i = 1, 2, \dots, m \right\},\,$$

then 
$$\lim_{k\to\infty}\mu_k(\Delta_\epsilon)=0.$$

We prove this by contradiction. Suppose the limit does not exist or is not 0, then there is a positive number C and a subsequence  $(\mu_{k_i})_{i\in\mathbb{N}}$  such that  $\mu_{k_i}(\Delta_{\epsilon}) > C$  for all i.

We define a linear functional  $\mathscr{L}(\mu) := \mathbb{E}_{\mu} f(\theta)$ , where  $f(\theta) = ||\theta - u||_2^2$  with  $u = \frac{\mathbf{e}}{m}$  the center of  $\Delta^{m-1}$ .

By definition, for a random variable following uniform distribution on  $\mathcal{D}$ ,  $\mathscr{L}\left(\Psi_{D}^{\mathbf{L}}(\mu)\right) = \mathbb{E}_{\mu}\left(\mathbb{E}_{d \sim D} f(T_{d}(\theta))\right)$ .

Consider that f is a strictly convex function, by Jensen's inequality,  $\mathbb{E}_{d\sim D}f(T_d(\theta))\geq f(\mathbb{E}_{d\sim D}T_d(\theta))=f(\theta)$ , with equality if and only if  $T_d(\theta)=\theta$  for all  $d\in \mathcal{D}$ , equivalently by the assumptions on matrix  $\mathbf{L}, \theta=\delta_h$  for some  $h\in \mathcal{H}$ . (This is because we assume  $\mathbf{L}$  have distinct columns, thus not all 2-by-2 cross-ratios are 1, for any pair of columns. however, after  $T_d$  all 2-by-2 cross-ratios are 1, indicating the existence of degeneration on every pair of columns. This can only happen when  $\theta=\delta_h$  for some  $h\in \mathcal{H}$ .)

Thus for any  $\theta \in \Delta_{\epsilon}$ ,  $\mathbb{E}_{d \sim D} f(T_d(\theta)) > f(\theta)$ . As  $\Delta_{\epsilon}$  is compact, there is a lower bound B > 0, such that

 $\mathbb{E}_{d\sim D} f(T_d(\theta)) - f(\theta) > B \text{ for all } \theta \in \Delta_{\epsilon}.$  Thus

$$\mathcal{L}(\mu_{k_i+1})$$

$$= \int_{\Delta^{m-1} \setminus \Delta_{\epsilon}} \mathbb{E}_{d \sim \mathcal{D}} f(T_d(\theta)) d\mu_{k_i} + \int_{\Delta_{\epsilon}} \mathbb{E}_{d \sim \mathcal{D}} f(T_d(\theta)) d\mu_{k_i}$$

$$> \int_{\Delta^{m-1} \setminus \Delta_{\epsilon}} f(\theta) d\mu_{k_i} + \int_{\Delta_{\epsilon}} f(\theta) d\mu_{k_i} + BC$$

$$= \mathcal{L}(\mu_{k_i}) + BC$$

for all  $i \in \mathbb{N}$  and simply  $\mathcal{L}(\mu_{k+1}) \geq \mathcal{L}(\mu_k)$  for general k.

Therefore,  $\mathcal{L}(\mu_k)$  is unbounded as  $k \to \infty$  since there is at least a BC > 0 increment at each  $k_i$ .

However, by definition, f is bounded by m since  $\sqrt{m}$  is the diameter of  $\Delta^{m-1}$  under 2-norm, thus  $\mathcal{L}(\mu) \leq m$ .

Such a contradiction shows that the opposite of our assumption,  $\lim_{k\to\infty}\mu_k(\Delta_\epsilon)=0$ , is valid.

Consider that  $\epsilon$  is arbitrary, and in Lemma C.1.1 we show that  $\mathbb{E}_{\mu_k}\theta(h)$  is invariant, thus  $\mu_k$  approaches  $\sum_{i\in\mathcal{H}}a_i\delta_i$  with  $a_i=\mathbb{E}_{\mu_0}\theta(i)$  in probability.

# C.2. Empirical Data for Stability: Perturbation on Prior

We sample 5 matrices of size  $3 \times 3$ , each of them are column-normalized, and their columns are sampled independently and uniformly on  $\Delta^2$ , listed below:

$$\mathbf{M}_{1} = \begin{pmatrix} 0.6559 & 0.5505 & 0.7310 \\ 0.1680 & 0.3359 & 0.0403 \\ 0.1760 & 0.1136 & 0.2287 \end{pmatrix}$$

$$\mathbf{M}_{2} = \begin{pmatrix} 0.2461 & 0.6600 & 0.4310 \\ 0.6785 & 0.0655 & 0.2325 \\ 0.0754 & 0.2746 & 0.3365 \end{pmatrix}$$

$$\mathbf{M}_{3} = \begin{pmatrix} 0.7286 & 0.1937 & 0.7620 \\ 0.0739 & 0.4786 & 0.1999 \\ 0.1974 & 0.3277 & 0.0382 \end{pmatrix}$$

$$\mathbf{M}_{4} = \begin{pmatrix} 0.4745 & 0.2024 & 0.5946 \\ 0.2898 & 0.7499 & 0.1313 \\ 0.2357 & 0.0477 & 0.2741 \end{pmatrix}$$

$$\mathbf{M}_{5} = \begin{pmatrix} 0.2207 & 0.5466 & 0.1605 \\ 0.3828 & 0.3807 & 0.5697 \\ 0.3965 & 0.0727 & 0.2698 \end{pmatrix}$$

$$(47)$$

And the 5 sampled priors are:

$$\theta_{1} = (0.3333, 0.3333, 0.3333)^{\top}$$

$$\theta_{2} = (0.1937, 0.4291, 0.3771)^{\top}$$

$$\theta_{3} = (0.4544, 0.0814, 0.4641)^{\top}$$

$$\theta_{4} = (0.5955, 0.2995, 0.1051)^{\top}$$

$$\theta_{5} = (0.4771, 0.0593, 0.4636)^{\top}$$
(48)

These names (including the rank 4 samples below) are overriding the previously defined identical symbols in this part and in the corresponding subsection in the main paper. In the  $4\times 4$  cases, we sample 3 matrices in the same way as in  $3\times 3$  case.

$$\mathbf{M}_{1}' = \begin{pmatrix} 0.3916 & 0.2306 & 0.0460 & 0.0404 \\ 0.1408 & 0.6350 & 0.2139 & 0.2310 \\ 0.2375 & 0.0275 & 0.1667 & 0.2412 \\ 0.2301 & 0.1068 & 0.5734 & 0.4874 \end{pmatrix}$$

$$\mathbf{M}_{2}' = \begin{pmatrix} 0.3744 & 0.6892 & 0.0112 & 0.3200 \\ 0.3204 & 0.2320 & 0.4498 & 0.3530 \\ 0.0291 & 0.0688 & 0.3865 & 0.0653 \\ 0.2761 & 0.0100 & 0.1526 & 0.2618 \end{pmatrix}$$

$$\mathbf{M}_{3}' = \begin{pmatrix} 0.2885 & 0.0873 & 0.2319 & 0.1009 \\ 0.0653 & 0.2239 & 0.0575 & 0.2584 \\ 0.5934 & 0.3276 & 0.2283 & 0.3925 \\ 0.0529 & 0.3612 & 0.4823 & 0.2482 \end{pmatrix} (49)$$

And 3 corresponding priors are sampled:

$$\theta'_1 = (0.2500, 0.2500, 0.2500, 0.2500)^{\top}$$
  

$$\theta'_2 = (0.1789, 0.3664, 0.2915, 0.1632)^{\top}$$
  

$$\theta'_3 = (0.4460, 0.4676, 0.0821, 0.0043)^{\top}$$
 (50)

The value we use to test the effectiveness of perturbed SCBI is called the successful rate, which is  $\mathbb{E}\left[\theta_{\infty}^{L}(h)\right] = \mathbb{E}_{\mu_{\infty}^{L}}\left[\theta(h)\right]$  where h is the true hypothesis that the teacher teaches (Definition 5.2). Successful rate is well defined, i.e. the limit exists, according to the convergences in probability (Theorem A.3 and Proposition C.1) with an  $\epsilon$  discussion based on them. To find the successful rate, we use Monte-Carlo method on  $10^4$  teaching sequences, and use Proposition C.1 to accelerate the simulation.

We can estimate an upper bound of the standard deviation (precision) of the empirical successful rate calculated based on Proposition C.1. The successful rate of a single teaching sequence is between 0 and 1, thus with a standard deviation smaller than 1. So the standard deviation of the empirical successful rate is bounded by  $(N)^{-1/2}$  where N is the number of sample sequences. Actually the precision is much

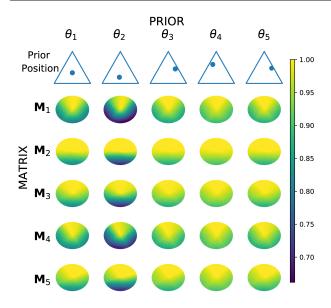


Figure 9. Successful rate of SCBI perturbed on prior. Each entry corresponds to a pair M and  $\theta_0^T$ . The first row shows for each prior  $\theta_0^T$ , the position it locates in  $\Delta^2$  and the range of  $\theta_0^L$  in the simulation. The 5 rows below are the zoom-in version of the shaded area in each case, whose color at each point represents the successful rate when  $\theta_0^L$  locates at that point.

smaller since the successful rate for a single sequence is much more stable.

Our first simulation is shown in Fig. 9, where we take  $\theta_0^L$  evenly on a series of concentric circles centered at  $\theta_0^T$ . There are 14 such circles with radius 0.005 to 0.07. On the *i*-th layer (smallest circle is the first layer) we take 6i many points evenly separated, the upper right figure in Fig. 10 shows how the points are taken in detail.

Thus we have 6 groups of points each distributed along a ray. We plot the successful rate versus the distance from the center along each ray in Fig. 10 for all the 25 combinations of  $\mathbf{M}$  and  $\theta_0^T$ .

To have a similar directional data for matrices of size  $4\times 4$ , we take a sample of 15 directions in  $\mathbb{R}^3$  (showing in Fig. 11 in spherical coordinates centered at  $(1/4,1/4,1/4,1/4,1/4)^{\top}$ , with  $(1,0,0,0)^{\top}$  as  $\phi=0$  axis and  $(0,1,0,0)^{\top}$  on the halfplane given by  $\theta=0$ ) and simulate the perturbations of  $\theta_0^L$  in  $\Delta^3$  along the 15 directions. On each ray, we take 20 evenly placed  $\theta_0^L$  with distance to the center  $\theta_0^T$  from 0.005 to 0.100. Then we plot the successful rate versus the distance in Fig. 11 for all 9 cases as before.

*Remark* 7. This part provides evidences of linear influence of the perturbation distance on the successful rate along a fixed direction.

Next we explore the global behavior of perturbations on prior. Here we sample for each combination of  $\mathbf{M}$  and  $\theta_0^T$  a set of 300 points for  $\theta_0^L$  evenly distributed in  $\Delta^3$ .

In Fig. 12, we plot the successful rate versus the value of  $\theta_0^L(h)$ , for all 25 situations.

We plot in Fig. 13 the distance to center as x-coordinates, for 9 situations with matrices of size  $4 \times 4$ .

In this part, we observe that there is a lower bound of the successful rate which depends linearly on the distance to center, with slope bounded by  $\frac{1}{\theta_n^2(h)}$  (Conjecture 5.3).

## C.3. Empirical Data for Stability: Perturbation on Matrix

Fig. 14 shows the behavior of perturbations on all sampled  $3\times 3$  matrices in Section 5. Perturbations are taken only along the relevant column / irrelevant column, since a perturbation on the target column is equivalent of the combination of a perturbation on other two columns (they have the same set of Cross-ratios, which determines the SCBI behavior). The cycle path in each plot is the equi-normalized-KL path, with any point on the path having the same normalized-KL to the target column as that of the original matrix  ${\bf T}$ .

These graphs should not be confused with the ones occur in the prior perturbed part, as we are plotting each column of the matrix here (the simplex is actually  $\mathcal{P}(\mathcal{D})$ ), while we were plotting the priors in previous discussion (the simplex is  $\mathcal{P}(\mathcal{H})$ ).

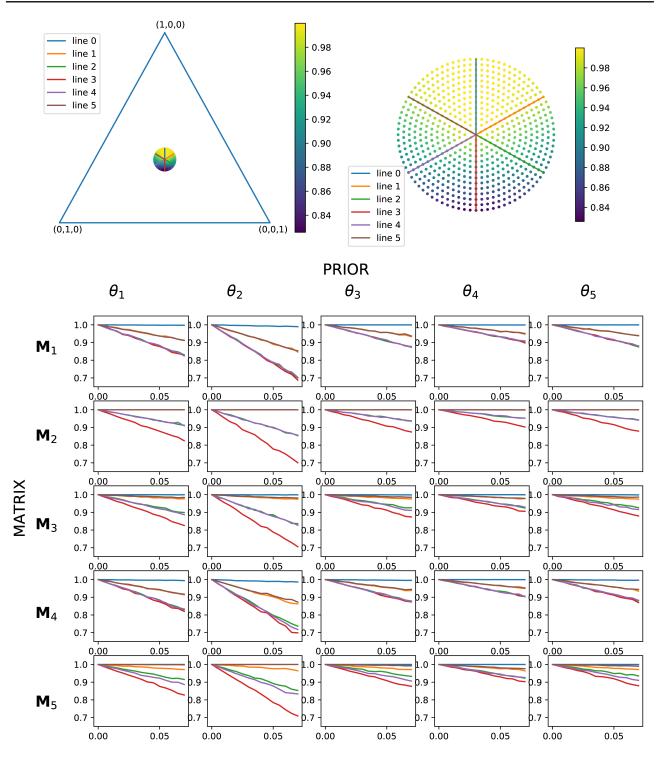


Figure 10. Upper-Left: the six rays at the center  $\theta_1$ . Upper-Right: zoom-in figure of the six rays in general. Lower: Successful rate versus distance to center along 6 rays. Fig. 5 in the main paper contains the Row 3 Column 1 picture of the lower one (with a different y-scale).

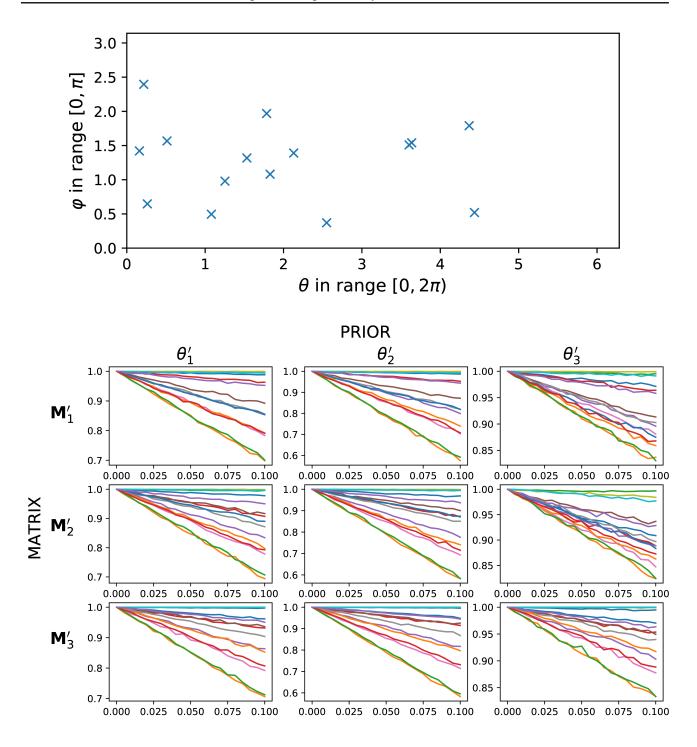


Figure 11. Upper: the sampled directions in spherical coordinates. Lower: Successful rate versus distance to center, along 15 rays, for all the 9 cases of matrices of size  $4 \times 4$ . The plot at Row 1 Column 1 appears in Fig. 5 of the main file.

## Acknowledgements

This material is based on research sponsored by the Air Force Research Laboratory and DARPA under agreement

number FA8750-17-2-0146 and the Army Research Office and DARPA under agreement HR00112020039. The U.S. Government is authorized to reproduce and distribute

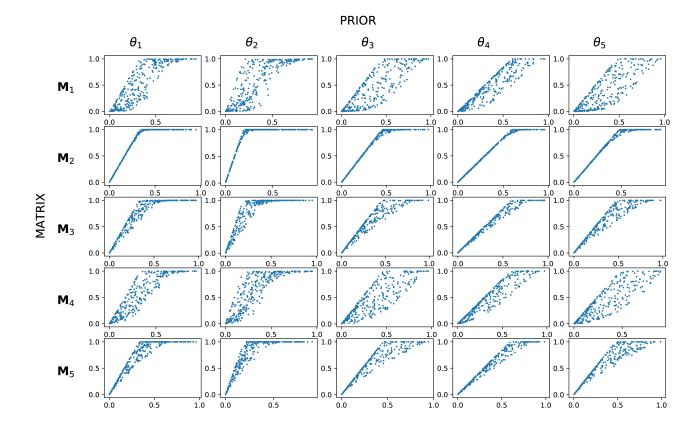


Figure 12. The 25 cases of  $3 \times 3$  matrices, with successful rate versus  $\theta_0^T(h)$  plotted. Plot at Row 3 Column 1 appears in Fig. 5 of the main file.

reprints for Governmental purposes notwithstanding any copyright notation thereon. This work was also supported by DoD grant 72531RTREP and NSF MRI 1828528 to PS.

## References

Bengio, Y., Louradour, J., Collobert, R., and Weston, J. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pp. 41–48. ACM, 2009.

Berger, J. O., Moreno, E., Pericchi, L. R., Bayarri, M. J., Bernardo, J. M., Cano, J. A., De la Horra, J., Martín, J., Ríos-Insúa, D., Betrò, B., et al. An overview of robust bayesian analysis. *Test*, 3(1):5–124, 1994.

Bonawitz, E., Shafto, P., Gweon, H., Goodman, N. D., Spelke, E., and Schulz, L. The double-edged sword of pedagogy: Instruction limits spontaneous exploration and discovery. *Cognition*, 120(3):322–330, 2011.

Bridgers, S., Jara-Ettinger, J., and Gweon, H. Young children consider the expected utility of others' learning to decide what to teachn. *Nature Human Behaviour*, in press.

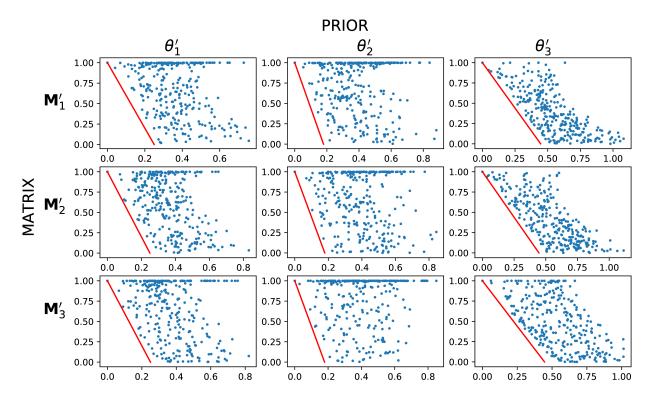


Figure 13. The 9 cases of  $4 \times 4$  matrices, with successful rate versus the distance to the center plotted. Plot at Row 1 Column 1 appears in Fig. 5 of the main file. Red line is the lower bound given in Conjecture 5.3.

Csibra, G. and Gergely, G. Natural pedagogy. *Trends in cognitive sciences*, 13(4):148–153, 2009.

Doliwa, T., Fan, G., Simon, H. U., and Zilles, S. Recursive teaching dimension, VC-dimension and sample compression. *Journal of Machine Learning Research*, 15(1): 3107–3131, 2014.

Doob, J. L. Application of the theory of martingales. *Le calcul des probabilites et ses applications*, pp. 23–27, 1949.

Eaves, B. S. and Shafto, P. Parameterizing developmental changes in epistemic trust. *Psychonomic Bulletin & Review*, pp. 1–30, 2016.

Eaves Jr, B. S., Feldman, N. H., Griffiths, T. L., and Shafto,P. Infant-directed speech is consistent with teaching.Psychological review, 123(6):758, 2016.

Elman, J. L. Learning and development in neural networks: The importance of starting small. *Cognition*, 48(1):71–99, 1993.

Fienberg, S. E. et al. An iterative procedure for estimation in contingency tables. *The Annals of Mathematical Statistics*, 41(3):907–917, 1970.

Fisac, J. F., Gates, M. A., Hamrick, J. B., Liu, C., Hadfield-Menell, D., Palaniappan, M., Malik, D., Sastry, S. S., Griffiths, T. L., and Dragan, A. D. Pragmatic-pedagogic value alignment. arXiv preprint arXiv:1707.06354, 2017.

Frank, M. C. and Goodman, N. D. Predicting pragmatic reasoning in language games. *Science*, 336(6084):998–998, 2012.

Ghahramani, Z. Probabilistic machine learning and artificial intelligence. *Nature*, 521(7553):452, 2015.

Goodman, N. D. and Stuhlmüller, A. Knowledge and implicature: Modeling language understanding as social cognition. *Topics in cognitive science*, 5(1):173–184, 2013.

Hadfield-Menell, D., Russell, S. J., Abbeel, P., and Dragan, A. Cooperative inverse reinforcement learning. In *Advances in neural information processing systems*, pp. 3909–3917, 2016.

- Hershkowitz, D., Rothblum, U. G., and Schneider, H. Classifications of nonnegative matrices using diagonal equivalence. *SIAM journal on Matrix Analysis and Applications*, 9(4):455–460, 1988.
- Ho, M. K., Littman, M., MacGlashan, J., Cushman, F., and Austerweil, J. L. Showing versus doing: Teaching by demonstration. In *Advances in Neural Information Processing Systems*, pp. 3027–3035, 2016.
- Jara-Ettinger, J., Gweon, H., Schulz, L. E., and Tenenbaum, J. B. The naïve utility calculus: Computational principles underlying commonsense psychology. *Trends in cognitive* sciences, 20(8):589–604, 2016.
- Kadane, J. B., Chuang, D. T., et al. Stable decision problems. *The Annals of Statistics*, 6(5):1095–1110, 1978.
- Miescke, K.-J. and Liese, F. *Statistical Decision Theory: Estimation, Testing, and Selection*. Springer, 2008. doi: https://doi.org/10.1007/978-0-387-73194-0.
- Murphy, K. P. *Machine learning: a probabilistic perspective*. MIT press, 2012.
- Schneider, M. H. Matrix scaling, entropy minimization, and conjugate duality. i. existence conditions. *Linear Algebra and its Applications*, 114:785–813, 1989.
- Shafto, P. and Goodman, N. Teaching games: Statistical sampling assumptions for learning in pedagogical situations. In *Proceedings of the 30th annual conference of the Cognitive Science Society*, pp. 1632–1637. Cognitive Science Society Austin, TX, 2008.
- Shafto, P., Goodman, N. D., and Frank, M. C. Learning from others: The consequences of psychological reasoning for human learning. *Perspectives on Psychological Science*, 7(4):341–351, 2012.
- Shafto, P., Goodman, N. D., and Griffiths, T. L. A rational account of pedagogical reasoning: Teaching by, and learning from, examples. *Cognitive Psychology*, 71:55–89, 2014.
- Skinner, B. F. Teaching machines. *Science*, 128(3330): 969–977, 1958.
- Tenenbaum, J. B., Kemp, C., Griffiths, T. L., and Goodman, N. D. How to grow a mind: Statistics, structure, and abstraction. *science*, 331(6022):1279–1285, 2011.
- Tomasello, M. *The cultural origins of human cognition*. Harvard University Press, Cambridge, MA, 1999.
- Wang, P., Paranamana, P., and Shafto, P. Generalizing the theory of cooperative inference. *AIStats*, 2019a.

- Wang, P., Wang, J., Paranamana, P., and Shafto, P. A mathematical theory of cooperative communication, 2019b.
- Yang, S. C., Yu, Y., Givchi, A., Wang, P., Vong, W. K., and Shafto, P. Optimal cooperative inference. In AISTATS, volume 84 of Proceedings of Machine Learning Research, pp. 376–385. PMLR, 2018.
- Zhu, X. Machine teaching for bayesian learners in the exponential family. In *Advances in Neural Information Processing Systems*, pp. 1905–1913, 2013.
- Zhu, X. Machine teaching: An inverse problem to machine learning and an approach toward optimal education. In *AAAI*, pp. 4083–4087, 2015.
- Zilles, S., Lange, S., Holte, R., and Zinkevich, M. Teaching dimensions based on cooperative learning. In *COLT*, pp. 135–146, 2008.

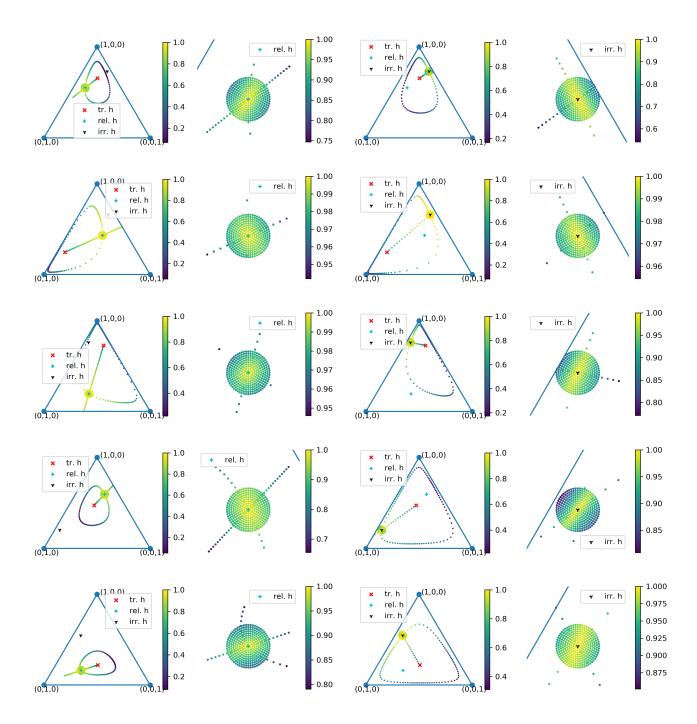


Figure 14. Perturbations on matrix **L**. First column: Perturbations on the irrelevant column of **L**. Second column: zoom-in of the first row. Third column: Perturbations on the relevant column of **L**. Last column: zoom-in of the third column. The scales of color in the zoomed figures are different from that of the original ones. Fig. 6 in the main paper is the third row here.