On the Performance of Social Learning

Felice Scala*, Marco Carpentiero*, Vincenzo Matta*, Ali H. Sayed[†]
*DIEM, University of Salerno, Fisciano, Italy

†EPFL, Lausanne, Switzerland

Abstract—Decentralized decision-making (a.k.a. social learning) deals with a group of agents, connected according to a graph, which cooperate to form their beliefs about some hypotheses of interest. Under mild technical conditions regarding the graph connectivity and the identifiability of the statistical model, the belief of each agent ultimately places all the probability mass on the true hypothesis when sufficient time to learn is granted. Less is known as regards the evaluation of the learning performance. One criterion that has been proposed is the rejection rate, i.e., the rate at which the belief about the wrong hypotheses converges to zero. We show in this work that this metric is not appropriate, since it leads to the paradoxical conclusion that the optimal Bayesian decision rule can be defeated by other rules. In contrast, we show that proper performance measures are the error probability and the error exponent, namely, the rate at which the error probability converges to zero. We compare different schemes in terms of these metrics, establishing useful connections between decentralized implementations and the optimal Bayesian system. Several interesting phenomena emerge. For example, we show that traditional social learning can be sensitive to the initial state and that the recently proposed doubly Non-Bayesian (NB²) learning scheme solves this issue.

Index Terms—Social learning, error probability, large deviations, rejection rate.

I. INTRODUCTION AND RELATED WORK

In the last few years, the theory of social learning (SL) [1]–[8] has been developed to characterize distributed decision-making systems. In these systems, spatially dispersed agents wish to accomplish a classification task to decide which hypothesis, chosen from a finite collection, is actually generating the observed data. To this end, the agents form their beliefs about each hypothesis (i.e., they assign probability scores to the hypotheses) and then opt for the hypothesis that maximizes their beliefs. The agents are allowed to cooperate only locally, i.e., with their neighbors. Moreover, for privacy or computational/memory constraints, they cannot share the raw data or the private decision models, but only their beliefs.

One fundamental result in SL theory is the following: given an infinite stream of data, under mild technical conditions the belief of each agent about the true hypothesis ultimately places the whole probability mass on the true hypothesis. It is also shown that the belief about any wrong hypothesis vanishes exponentially fast, with a certain *rejection rate* [5], [6]. Based on this type of result, there are works in the literature that compare different SL strategies in terms of the rejection rates [9], [10]. For example, in [9] a strategy is proposed that can

The work of V. Matta was partially supported by the European Union under the Italian National Recovery and Resilience Plan (NRRP) of NextGenerationEU, partnership on "Telecommunications of the Future" (PE00000001 - program "RESTART").

exhibit a larger rejection rate than traditional SL. However, traditional SL is known to attain (in the precise sense specified later) the optimal performance guaranteed by the centralized Bayesian posterior. This would imply that the scheme proposed in [9] would outperform the optimal centralized system. raising the fundamental question of whether the rejection rate is an appropriate performance measure. One contribution of this work is to show that it is actually not. Motivated by this finding, we then envisage that the error probability is one meaningful performance criterion. However, due to the lack of analytical relations to compute the error probability, we replace it by the error exponent, which quantifies the rate at which the error probability vanishes as the size of the data streams increases. Recent results [11] have shown that traditional social learning attains the optimal error exponent of the centralized Bayesian posterior only when the combination matrix is doubly stochastic. When it is simply left stochastic, a different strategy is available [11] that guarantees asymptotic optimality. As a second contribution, we show that asymptotic optimality in terms of error exponents can hide factors that matter in the error probability. In particular, we focus on the effect of the initial belief vectors. This effect is commonly overlooked. This is because, in the literature on decentralized optimization and learning, the initial states are asymptotically washed out and have no effect on the mean-square-error and related performance metrics. We ascertain that this is no longer the case in the decision-making setting. Specifically, we show that the impact of the initial beliefs is washed out in terms of the error exponents, but *not* in terms of the error probabilities. We also show that the NB² strategy proposed in [11] is able to remediate this problem.

II. BACKGROUND AND NOTATION

Let $\Theta = \{\theta_1, \dots, \theta_H\}$ be a set of H hypotheses. In social learning, K agents cooperate to discover the true state of nature, θ_0 , assumed included within the set Θ . Each agent $k=1,\dots,K$, at time $t=1,2,\dots$, has access to a local stream of data samples $x_{k,t} \in \mathcal{X}_k$ (we use bold notation for random quantities). The data are independent over time, and may be dependent across the agents. Ideally, to solve the decision-making problem optimally, the agents should compute the true posterior probabilities for any $\theta \in \Theta$:

$$\boldsymbol{\mu}_t^{\star}(\theta) \propto \pi(\theta) \prod_{\tau=1}^t \ell_{\text{tot}}(\boldsymbol{x}_{1,\tau}, \dots, \boldsymbol{x}_{K,\tau} | \theta),$$
 (1)

where $\pi(\theta)$ is the prior probability of θ , $\ell_{tot}(\cdot|\theta)$ is the *joint* likelihood model, and the symbol ∞ hides the proportion-

ality constant needed to make $\mu_t^{\star} = [\mu_t^{\star}(\theta_1), \dots, \mu_t^{\star}(\theta_H)]$ a probability vector. The true distribution of $[x_{1,t}, \dots, x_{K,t}]$ is $\ell_{\text{tot}}(\cdot|\theta_0)$. Usually, the true posterior given by (1) cannot be implemented in a decentralized setting, especially because the agents cannot share their raw data, they are allowed to communicate only with their neighbors, and are unable to build the joint likelihood model [1].

To perform social learning, each agent k employs instead a marginal likelihood model $\ell_k(\cdot|\theta)$. The agents' opinions about the hypotheses are represented by a belief vector, i.e., a vector of probability scores assigned to each $\theta \in \Theta$ at each time t, defined as $\mu_{k,t} = [\mu_{k,t}(\theta_1), \mu_{k,t}(\theta_2), \dots, \mu_{k,t}(\theta_H)]$. Then, agent k chooses the hypothesis that maximizes its belief, which, for the case of the true posterior, would correspond to the optimal maximum a posteriori probability (MAP) rule [12].

A. Traditional Social Learning

In traditional social learning, the belief vectors across the agents are updated iteratively by the following two steps:

$$\psi_{k,t}(\theta) \propto \mu_{k,t-1}(\theta) \, \ell_k(\boldsymbol{x}_{k,t}|\theta),$$
 (2a)

$$\mu_{k,t}(\theta) \propto \prod_{j=1}^{K} \left[\psi_{j,t}(\theta)\right]^{a_{jk}},$$
 (2b)

The beliefs are initialized with some deterministic values $\mu_{k,0}(\theta)$. The first step (2a) performs a local Bayesian update by using the fresh data sample $\boldsymbol{x}_{k,t}$, yielding an intermediate belief $\psi_{k,t}(\theta)$. The second step (2b) pools the beliefs of neighboring agents. Specifically, the beliefs are (geometrically) weighted by the nonnegative coefficients a_{jk} collected into the combination matrix $A = [a_{jk}] \in \mathbb{R}^{K \times K}$. The weights a_{jk} are nonzero only when agent k can receive information from agent k. We adopt the following standard assumption.

Assumption 1 (Combination matrix). The combination matrix A is i) left stochastic, which means that $a_{jk} \geq 0$ and $\sum_{j=1}^{K} a_{jk} = 1$; and ii) primitive, which means that it is irreducible and has a single eigenvalue (equal to 1) on the spectral circle [12].

The typical question is whether the SL strategy (2a)-(2b) allows *all* the agents in the network to detect the true hypothesis. This is possible if at least one agent is able to correctly detect the true hypothesis [1], [2].

Assumption 2 (Global identifiability). For each $\theta \neq \theta_0$, there exists at least one agent k such that $D(\ell_{k,\theta_0}||\ell_{k,\theta}) > 0$, where $D(\ell_{k,\theta_0}||\ell_{k,\theta})$ denotes the Kullback-Leibler (KL) divergence between $\ell_k(\cdot|\theta_0)$ and $\ell_k(\cdot|\theta)$. Moreover, the KL divergence between any pair of hypotheses is assumed finite.

Under Assumptions 1 and 2, the social learning strategy in (2a)-(2b) achieves truth learning in the sense that, for all k,

$$\mu_{k,t}(\theta_0) \xrightarrow[t \to \infty]{\text{a.s.}} 1,$$
 (3)

where $\xrightarrow[t\to\infty]{\text{a.s.}}$ denotes almost-sure convergence [1].

III. EXAMINING THE REJECTION RATE

Let $\mu_{k,t}$ be a belief vector obtained from some construction that can also be different from (2a)-(2b). Assume that, for all $\theta \neq \theta_0$,

$$\frac{1}{t}\log\frac{1}{\boldsymbol{\mu}_{k,t}(\theta)} \xrightarrow[t \to \infty]{\text{a.s.}} \rho(\theta, \theta_0) > 0, \tag{4}$$

then $\rho(\theta, \theta_0)$ is called the *rejection rate* (of θ when θ_0 is true). Note that (4) implies that $\mu_{k,t}(\theta_0)$ converges to 1 (exponentially) with probability 1.

While it is always possible to compare two strategies in terms of their rejection rates, this does not necessarily leads to a meaningful performance comparison since, as the following lemma indicates, it is always possible to modify either strategies to get an arbitrarily better rejection rate.

Lemma 1 (Arbitrariness of the rejection rate). Let $\mu_{k,t}$ be a belief vector achieving rejection rates $\rho(\theta, \theta_0)$, for $\theta \neq \theta_0$. Then, it is always possible to construct another belief vector $\tilde{\mu}_{k,t}$ that achieves rejection rates $\tilde{\rho}(\theta, \theta_0) > \rho(\theta, \theta_0)$.

Proof: Let, for an arbitrary b > 1,

$$\tilde{\boldsymbol{\mu}}_{k,t}(\theta) \propto \boldsymbol{\mu}_{k,t}^b(\theta).$$
 (5)

Observe that $\mu_{k,t}(\theta_0)$ converges almost surely to 1 by assumption and, hence, so does $\tilde{\mu}_{k,t}(\theta_0)$. Then, to compute the rejection rate we can compute the limit of $(1/t)\log\frac{\tilde{\mu}_{k,t}(\theta_0)}{\tilde{\mu}_{k,t}(\theta)}$. We have, for any $\theta \neq \theta_0$:

$$\frac{1}{t}\log\frac{\tilde{\boldsymbol{\mu}}_{k,t}(\theta_0)}{\tilde{\boldsymbol{\mu}}_{k,t}(\theta)} = \frac{b}{t}\log\frac{\boldsymbol{\mu}_{k,t}(\theta_0)}{\boldsymbol{\mu}_{k,t}(\theta)} \xrightarrow[t \to \infty]{\text{a.s.}} b\rho(\theta,\theta_0) \triangleq \tilde{\rho}(\theta,\theta_0),$$
(6)

where $\tilde{\rho}(\theta, \theta_0) > \rho(\theta, \theta_0)$ since b > 1.

Lemma 1 reveals that, given an arbitrary SL strategy, it is always possible to modify that strategy by adding (5) to it and attain any desired rejection rate. This suggests that the rejection rate is not an effective metric to evaluate decision-making strategies.

Example 1 (One-bit quantized observations). Consider the following decentralized binary classification problem. Let $\boldsymbol{x}_{k,t}$ be unit-variance Gaussian random variables, independent and identically distributed (iid) across time and space, with means m_{θ_1} and m_{θ_2} under the two hypotheses. Let, for $\gamma \in \mathbb{R}$, $q_{k,t} = \mathbb{I}[\boldsymbol{x}_{k,t} > \gamma]$, where \mathbb{I} is the indicator function, which is equal to 1 if the condition defined by its argument is true, and is 0 otherwise. That is, $q_{k,t}$ is a one-bit quantized version of $\boldsymbol{x}_{k,t}$. Accordingly, the statistical distribution of $q_{k,t}$ is uniquely determined by the probability

$$\mathbb{P}[\boldsymbol{q}_{k,t} = 1|\theta] = Q(\gamma - m_{\theta}), \qquad \theta \in \{\theta_1, \theta_2\}, \tag{7}$$

where Q is the complementary cumulative distribution function of the standard normal. Since $q_{k,t}$ is a transformation of the original data $x_{k,t}$, it is obvious that any decision strategy based on $q_{k,t}$ cannot outperform the optimal decision strategy based on $x_{k,t}$. Nevertheless, according to Lemma 1, we can design a learning scheme based on $q_{k,t}$ that achieves

an arbitrarily large rejection rate. We now elaborate on this aspect. Denote by $\boldsymbol{\mu}_{k,t}^{(q)}(\theta)$ the belief of agent k at time t computed with the traditional social learning scheme (2a)-(2b) applied to the quantized observations. Accordingly, the marginal likelihood of each agent k is $\ell^{(q)}(y|\theta) = \mathbb{P}[q_{k,t} = y|\theta]$, for $y \in \{0,1\}$. Under Assumptions 1 and 2, and with nonzero initial beliefs, one can show that, for all $\theta \neq \theta_0$ [1]

$$\frac{1}{t}\log\frac{1}{\boldsymbol{\mu}_{k,t}^{(q)}(\theta)} \xrightarrow[t \to \infty]{\text{a.s.}} D\left(\ell_{\theta_0}^{(q)}||\ell_{\theta}^{(q)}\right), \tag{8}$$

The KL divergence on the RHS represents the rejection rate of the SL scheme based on the quantized variables. Likewise, it can be shown that the optimal (centralized) Bayesian strategy that computes the true posterior (1) based on the unquantized observations $x_{k,t}$ would reach a rejection rate $K \cdot D(\ell_{\theta_0} || \ell_{\theta})$ (where we suppressed the subscript k in the likelihoods due to the identical distribution across the agents) [1]. Observe that, in view of the data processing inequality for the KL divergence, we have $D(\ell_{\theta_0}||\ell_{\theta}) \geq D(\ell_{\theta_0}^{(q)}||\ell_{\theta}^{(q)})$ [13]. On the other hand, by applying Lemma 1, we can modify the scheme based on the quantized variables by constructing a belief $\tilde{\mu}_{k,t}^{(q)}(\theta) \propto [\mu_{k,t}^{(q)}(\theta)]^b$, with b chosen to attain a rejection rate $b \cdot D(\ell_{\theta_0}^{(q)}||\ell_{\theta}^{(q)}) \ge K \cdot D(\ell_{\theta_0}||\ell_{\theta})$, i.e., exceeding the rejection rate achieved by μ_t^{\star} . As a result, if the rejection rate is taken as a performance proxy, we reach the conclusion that a decentralized scheme based on quantized variables can outperform the optimal centralized scheme based on the original unquantized variables, which appears to be a suspicious result.

IV. ERROR PROBABILITY PERFORMANCE

The true posterior $\mu_t^{\star}(\theta)$ is the best information available regarding the characterization of the link between the data and the hypotheses. As a result, it allows us to make the best inference possible.

Assume now that we want to judge the goodness of a decision strategy in terms of the rejection rate. Let $\rho^*(\theta,\theta_0)$ be the rejection rate achieved by the true posterior. In view of Lemma 1, we can always find another strategy that would guarantee a rejection rate larger than $\rho^*(\theta,\theta_0)$. For example, with reference to Example 1, in place of the optimal posterior of the data $x_{k,t}$ we could use a strategy based on one-bit data and modify it so as to increase arbitrarily the rejection rate. This would imply that there are (actually, infinitely many) strategies that converge to the ideal belief vector faster than the optimal posterior. This suggests that the rejection rate is not appropriate to measure the goodness of an inferential strategy.

The optimal posterior tells us what is the best level of confidence we must place on a given decision. If the actual probability that θ is true given the observed data realization is equal to 0.9, we know from decision theory that 90% is the best level of confidence we must place on our decision. We should be neither more nor less confident. Using another strategy that increases the rejection rate by placing more mass on the true hypothesis is not a real gain, since we are mistakenly increasing our level of confidence. In other words,

the belief about the true hypothesis should converge to 1 as fast as dictated by the optimal posterior, but not faster.

One possible misunderstanding in the interpretation of the rejection rate arises from overlooking the randomness in the considered problem. Since the belief is random, whatever the speed of convergence (i.e., the rejection rate) is, there will always be a nonzero probability of error, which cannot be ignored when comparing different strategies. To gain insights, refer back to Example 1. The MAP rule based on (1) (i.e., on the unquantized data) minimizes the error probability. Consider next the belief $\mu_{k,t}^{(q)}$ obtained by applying traditional SL to the *quantized* data, and its modified version $\tilde{\mu}_{k+1}^{(q)}$ which achieves a rejection rate larger than the optimal MAP rule. Observe that the error events $\{\tilde{\boldsymbol{\mu}}_{k,t}^{(\mathbf{q})}(\theta_0) \leq \tilde{\boldsymbol{\mu}}_{k,t}^{(\mathbf{q})}(\theta)\}$ are equivalent to the events $\{\boldsymbol{\mu}_{k,t}^{(\mathbf{q})}(\theta_0) \leq \boldsymbol{\mu}_{k,t}^{(\mathbf{q})}(\theta)\}$, which means that $\pmb{\mu}_{k,t}^{(\mathrm{q})}$ and $\tilde{\pmb{\mu}}_{k,t}^{(\mathrm{q})}$ make the same mistakes. Moreover, they make more mistakes than the optimal posterior μ_t^{\star} . Thus, the higher rejection rate is not rewarding. This can be explained in terms of the randomness in the data. Indeed, consider the modified belief $ilde{\mu}_{k,t}^{(\mathrm{q})}$ and the optimal belief μ_t^{\star} at a given time instant t. For some data realizations, $\tilde{\mu}_{k,t}^{(q)}$ leads to correct decisions and, in view of its higher rejection rate, discards the wrong hypotheses with more confidence than μ_t^{\star} . However, this is not sufficient to outperform the MAP rule, since $ilde{oldsymbol{\mu}}_{k,t}^{(\mathrm{q})}$ is not optimal and, hence, there are more realizations where it decides wrongly while μ_t^* does not.

A. Error Probability and Error Exponent

From the previous discussion we conclude that a meaningful criterion to evaluate the performance of social learning strategies is the error probability achieved by each individual agent. Specifically, recalling that the decision is made by taking the hypothesis that maximizes the belief, the error probability of agent k given that the true hypothesis is θ_0 is given by

$$p_{k,t}(\theta_0) = \mathbb{P}_{\theta_0}[\boldsymbol{\mu}_{k,t}(\theta_0) \le \boldsymbol{\mu}_{k,t}(\theta) \text{ for at least one } \theta \ne \theta_0],$$
(9)

where the notation \mathbb{P}_{θ_0} signifies that the true hypothesis is θ_0 . Then, the total error probability of agent k is obtained by averaging over the prior distribution π , yielding

$$p_{k,t} = \sum_{\theta_0 \in \Theta} \pi(\theta_0) p_{k,t}(\theta_0). \tag{10}$$

Unfortunately, closed-form relations to express this probability are seldom available. However, suitable asymptotic (as $t \to \infty$) performance measures related to the error probability can be obtained by exploiting the asymptotic properties of the log belief ratios $\log(\mu_{k,t}(\theta_0)/\mu_{k,t}(\theta))$ [1]. One of the performance measures particularly suited to decision problems is the *error exponent*, which quantifies the rate at which the *error probability* converges (exponentially fast) to 0 as $t \to \infty$. This exponent is given by

$$\lim_{t \to \infty} \frac{1}{t} \log p_{k,t} = E,\tag{11}$$

which can be evaluated by resorting to the *large deviations* theory [14], [15]. Different from the rejection rate, the error exponent refers to convergence rate of the *error probability*.

It is shown in [11] that, for statistically independent agents, traditional social learning attains the same error exponent as the optimal centralized Bayesian solution when the combination matrix is doubly stochastic. When it is only left stochastic, it is necessary to modify traditional social learning into the socialled non-Bayesian-non-Bayesian (NB²) strategy

$$\psi_{k,t}^{(\mathrm{NB}^2)}(\theta) \propto \mu_{k,t-1}^{(\mathrm{NB}^2)}(\theta) \, \ell_k^{1/v_k}(\boldsymbol{x}_{k,t}|\theta),$$
 (12a)

$$\boldsymbol{\mu}_{k,t}^{(\mathrm{NB}^2)}(\boldsymbol{\theta}) \propto \prod_{j=1}^K \left[\boldsymbol{\psi}_{j,t}^{(\mathrm{NB}^2)}(\boldsymbol{\theta}) \right]^{a_{jk}}, \tag{12b}$$

which replaces the Bayesian update (2a) with a non-Bayesian update where the likelihood is raised to the k-th entry v_k of the Perron vector associated with the combination matrix A.

B. Limitations of the Error Exponent

As already indicated, the theory of large deviations provides appropriate tools to compute the error exponent. However, expressions like (11) are equivalent to $p_{k,t} = \exp{\{-E\,t + o(t)\}}$, where by o(t) we define a quantity such that $\lim_{t\to\infty} o(t)/t = 0$. This means that o(t) can be a constant or even a term that diverges to ∞ as $t\to\infty$. Accordingly, we may have strategies that share the same error exponent but have different error probabilities, which implies that the error exponent cannot be used to approximate error probabilities.

Example 2 (Role of the initial beliefs). Let $x_{k,t}$ follow the same Gaussian model described in Example 1. Assume that the agents run a traditional SL algorithm. We want to show that the error probability of traditional SL may suffer when the initial beliefs are non-uniform across the hypotheses. To magnify this specific aspect, we consider the favorable case where all agents have the same initial beliefs $\mu_{k,0}=\pi$ and the network is fully connected. We will see that even in this favorable case the uneven initial belief assignment matters.

By unfolding the recursion in (2a)-(2b), we can write, over a fully connected network with $a_{jk}=1/K$ for all j and k,

$$\log \frac{\boldsymbol{\mu}_{k,t}(\theta_1)}{\boldsymbol{\mu}_{k,t}(\theta_2)} = \xi + \frac{1}{K} \sum_{\tau=1}^{t} \sum_{k=1}^{K} \log \frac{\ell_k(\boldsymbol{x}_{k,\tau}|\theta_1)}{\ell_k(\boldsymbol{x}_{k,\tau}|\theta_2)},$$
(13)

where we defined $\xi \triangleq \log(\pi(\theta_1)/\pi(\theta_2))$. A similar expression is obtained for the true posterior in (1), namely,

$$\log \frac{\boldsymbol{\mu}_t^{\star}(\theta_1)}{\boldsymbol{\mu}_t^{\star}(\theta_2)} = \xi + \sum_{\tau=1}^t \sum_{k=1}^K \log \frac{\ell_k(\boldsymbol{x}_{k,\tau}|\theta_1)}{\ell_k(\boldsymbol{x}_{k,\tau}|\theta_2)}.$$
 (14)

It is apparent that (13) and (14) are basically the same, with the only difference being the factor 1/K that scales the sum of log likelihoods in (13). Moreover, it is also apparent that the term due to the initial beliefs is independent of t and becomes immaterial as $t \to \infty$, which would also suggest that the scaling factor 1/K becomes asymptotically immaterial. However, there is a subtle effect here that needs to be carefully

examined. To this end, let us evaluate the error probability. This is doable in this example, due to the Gaussian model.

Assume that the true hypothesis is θ_1 . Using (13), we can compute the error probability in (9) as

$$p_{k,t}(\theta_1) = \mathbb{P}_{\theta_1} \left[\frac{1}{K} \sum_{\tau=1}^{t} \sum_{k=1}^{K} \log \frac{\ell_k(\boldsymbol{x}_{k,\tau} | \theta_2)}{\ell_k(\boldsymbol{x}_{k,\tau} | \theta_1)} \ge \xi \right].$$
 (15)

It is readily verified that, when θ_1 is true, $\log \frac{\ell_k(\boldsymbol{x}_{k,\tau}|\theta_2)}{\ell_k(\boldsymbol{x}_{k,\tau}|\theta_1)}$ is Gaussian with mean $-\Delta$ and variance 2Δ , where $\Delta=(m_2-m_1)^2/2$ is the KL divergence between the two hypotheses [16]. Using this result in (15), we get

$$p_{k,t}(\theta_1) = Q\left(\frac{\xi + t\Delta}{\sqrt{2t\Delta/K}}\right) \approx \frac{1}{2}e^{-\frac{K\Delta t}{4} - \frac{K\xi}{2} - \frac{K\xi^2}{4\Delta t}}, \quad (16)$$

where we have approximated the Q-function as $Q(x) \approx 1/2 \exp\{-x^2/2\}$. Following similar steps, for the optimal Bayesian system we would get from (14)

$$p_t^{\star}(\theta_1) = Q\left(\frac{\xi + Kt\Delta}{\sqrt{2Kt\Delta}}\right) \approx \frac{1}{2}e^{-\frac{K\Delta t}{4} - \frac{\xi}{2} - \frac{\xi^2}{4K\Delta t}}.$$
 (17)

By symmetry, it is readily seen that the error probability when θ_2 is true correspond to replacing ξ with $-\xi$ in (16) and (17). Accordingly, from (10) we can write the error probabilities for the decentralized and the optimal system as

$$p_{k,t} \approx \frac{1}{2} e^{-\frac{K\Delta t}{4} - \frac{K\xi^2}{4\Delta t}} \left(\pi(\theta_1) e^{-\frac{K\xi}{2}} + (1 - \pi(\theta_1)) e^{\frac{K\xi}{2}} \right),$$
 (18)

$$p_t^{\star} \approx \frac{1}{2} e^{-\frac{K\Delta t}{4} - \frac{\xi^2}{4K\Delta t}} \left(\pi(\theta_1) e^{-\frac{\xi}{2}} + (1 - \pi(\theta_1)) e^{\frac{\xi}{2}} \right).$$
 (19)

Comparing (18) against (19), we see that, as $t \to \infty$, the two expressions are equivalent at the leading exponential order:

$$\lim_{t \to \infty} \frac{1}{t} \log p_{k,t}(t) = \lim_{t \to \infty} \frac{1}{t} \log p_t^* = \frac{K\Delta}{4}.$$
 (20)

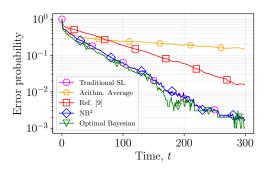
However, we also see that the term due to the initial belief is reduced by a factor K in the optimal Bayes scheme, which yields the following limit:

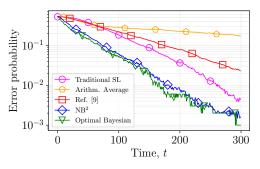
$$\lim_{t \to \infty} \frac{p_{k,t}}{p_t^*} \approx \frac{\pi(\theta_1)e^{-\frac{K\xi}{2}} + (1 - \pi(\theta_1))e^{\frac{K\xi}{2}}}{\pi(\theta_1)e^{-\frac{\xi}{2}} + (1 - \pi(\theta_1))e^{\frac{\xi}{2}}},$$
 (21)

which, using the definition of ξ , after some straightforward algebraic manipulations gives

$$\lim_{t \to \infty} \frac{p_{k,t}}{p_t^*} \approx \cosh\left(\frac{(K-1)\,\xi}{2}\right) > 0,\tag{22}$$

which is strictly positive for all K>1, revealing that, even asymptotically, the error probability of the decentralized system is larger than the error probability of the optimal system. This fact does not contradict the large deviation estimate (which shows that the two strategies share the same error exponent), since it is well known that the error exponent gives the leading-order exponential term, while neglecting possible higher-order corrections at the exponent. Notably, it can be shown that the NB² strategy (12a)-(12b) is able to compensate for the uneven initial belief assignment, since the factor v_k to which the likelihood is raised (in the considered case $v_k=1/K$) would actually turn (13) into (14).





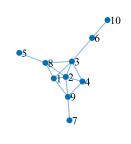


Fig. 1: Error probability $(1/K) \sum_{k=1}^{K} p_{k,t}$, where $p_{k,t}$ is defined by (10). Left: Uniform initial beliefs. Middle: Uneven initial beliefs, with $\pi(\theta_1) = 0.2$, $\pi(\theta_2) = 0.3$, $\pi(\theta_3) = 0.5$ for all agents. Right Network topology.

V. ILLUSTRATIVE EXAMPLES

In this section we present the results of some numerical experiments conducted in the following setup. We consider the network topology shown in Fig. 1. On top of this topology, we build a doubly stochastic combination matrix by applying the Metropolis policy [12]. The decision problem is a ternary classification problem. The observations are iid over time and across the agents. Each observation $x_{k,t}$ is a five-dimensional vector whose entries are iid unit-variance Gaussian variables with means 0, 0.05, and 0.1 under hypotheses θ_1 , θ_2 , and θ_3 .

We will compare the following strategies: traditional SL (2a)-(2b); the strategy that, in step (2b), replaces the weighted geometric average of the intermediate beliefs with their arithmetic average [1], [3]; the SL algorithm proposed in [9]; the NB² strategy (12a)-(12b); and the optimal Bayesian classifier.

The performance indices used to compare these strategies are the error probability and the error exponent. Before examining the results, let us summarize the guarantees that we have from the theoretical analysis. We know that: i) NB² attains the optimal error exponent (this is also true for traditional social learning when the combination matrix is doubly stochastic). As a result, no other strategy can outperform NB² in terms of error exponent; ii) the gain in error exponent can neglect sub-exponential corrections that matters in terms of error probability. For example traditional SL suffers from uneven initial beliefs, but NB² is able to compensate for this effect.

In the left panel of Fig. 1, we compare the five strategies when the initial beliefs are uniform. As guaranteed by the theoretical results, traditional SL and NB² guarantee the same error exponent as the optimal Bayesian classifier. Since no strategy can beat the optimal Bayes system, no strategy can have a better exponent. Accordingly, arithmetic averaging and the SL algorithm in [9] feature a worse (i.e., smaller) error exponent. In this particular example, arithmetic averaging is outperformed by the SL algorithm in [9].

In terms of error probabilities, the better exponent allows traditional SL and NB^2 to outperform the other two strategies. Moreover, since the initial beliefs are uniform, the error probabilities of traditional SL and NB^2 are also equal.

The case of uneven initial beliefs is addressed in the middle panel of Fig. 1. Comparing with the left panel, we see that traditional SL suffers from the uneven initial beliefs, as discussed in Example 2. We also see that NB² remediates this issue and gets closer to the optimal error probability. The comments and relative ordering pertaining to arithmetic averaging SL and to the algorithm in [9] are unchanged.

REFERENCES

- [1] V. Matta, V. Bordignon, and A. H. Sayed, *Social Learning: Opinion Formation and Decision-Making Over Graphs*. Now Publishers, 2025.
- [2] V. Bordignon, V. Matta, and A. H. Sayed, "Socially intelligent networks: A framework for decision making over graphs," *IEEE Sig. Process. Mag.*, vol. 41, no. 4, pp. 20–39, Jul. 2024.
- [3] X. Zhao and A. H. Sayed, "Learning over social networks via diffusion adaptation," in *Proc. Asilomar Conference on Signals, Systems, and Computers*, Pacific Grove, CA, USA, 2012, pp. 709–713.
- [4] A. Jadbabaie, P. Molavi, A. Sandroni, and A. Tahbaz-Salehi, "Non-Bayesian social learning," *Games and Economic Behavior*, vol. 76, no. 1, pp. 210–225, Sep. 2012.
- [5] A. Lalitha, T. Javidi, and A. D. Sarwate, "Social Learning and Distributed Hypothesis Testing," *IEEE Trans. Inf. Theory*, vol. 64, no. 9, pp. 6161–6179, Sep. 2018.
- [6] A. Nedić, A. Olshevsky, and C. A. Uribe, "Fast Convergence Rates for Distributed Non-Bayesian Learning," *IEEE Trans. Automat. Control*, vol. 62, no. 11, pp. 5538–5553, Nov. 2017.
- [7] C. Chamley, Rational Herds: Economic Models of Social Learning. Cambridge University Press, 2004.
- [8] V. Krishnamurthy and H. V. Poor, "Social learning and Bayesian games in multiagent signal processing: How do local and global decision makers interact?" *IEEE Sig. Process. Mag.*, vol. 30, no. 3, pp. 43–57, May 2013.
- [9] A. Mitra, J. A. Richards, and S. Sundaram, "A new approach to distributed hypothesis testing and non-Bayesian learning: Improved learning rate and Byzantine resilience," *IEEE Trans. Automat. Control*, vol. 66, no. 9, pp. 4084–4100, Sep. 2021.
- [10] M. Kayaalp, Y. İnan, E. Telatar, and A. H. Sayed, "On the arithmetic and geometric fusion of beliefs for distributed inference," *IEEE Trans. Automat. Control*, vol. 69, no. 4, pp. 2265–2280, Apr. 2024.
- [11] V. Bordignon, M. Kayaalp, V. Matta, and A. H. Sayed, "Social learning with non-Bayesian local updates," in *Proc. EUSIPCO*, Helsinki, Finland, Sep. 2023, pp. 1–5.
- [12] A. H. Sayed, Inference and Learning from Data, 3 vols., Cambridge University Press, 2022.
- [13] Y. Polyanskiy and Y. Wu, Information Theory: From Coding to Learning. Cambridge University Press, 2023.
- [14] A. Dembo and O. Zeitouni, Large Deviations Techniques and Applications. Springer, 2009.
- [15] F. den Hollander, Large Deviations. AMS, 2000.
- [16] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. John Wiley & Sons, 1991.