

Consistency in Models for Distributed Learning under Communication Constraints

Joel B. Predd, *Member, IEEE*, Sanjeev R. Kulkarni, *Fellow, IEEE*,
and H. Vincent Poor, *Fellow, IEEE*

Abstract

Motivated by sensor networks and other distributed settings, several models for distributed learning are presented. The models differ from classical works in statistical pattern recognition by allocating observations of an independent and identically distributed (i.i.d.) sampling process amongst members of a network of simple learning agents. The agents are limited in their ability to communicate to a central fusion center and thus, the amount of information available for use in classification or regression is constrained. For several basic communication models in both the binary classification and regression frameworks, we question the existence of agent decision rules and fusion rules that result in a universally consistent ensemble; the answers to this question present new issues to consider with regard to universal consistency. This paper addresses the issue of whether or not the guarantees provided by Stone's Theorem in centralized environments hold in distributed settings.

Index Terms

Classification, regression, consistency, distributed learning, nonparametric, sensor networks, statistical pattern recognition

This paper was presented in part at the 17th Annual Conference on Learning Theory (COLT), Banff, Canada, July 1-4, 2004 [31] and in part at the 42nd Annual Allerton Conference on Communication, Control, and Computing, Monticello, IL, Sept 29-Oct 1, 2004 [32].

This research was supported in part by the Army Research Office under grant DAAD19-00-1-0466, in part by Draper Laboratory under grant IR&D 6002, in part by the National Science Foundation under grant CCR-0312413, and in part by the Office of Naval Research under Grant No. N00014-03-1-0102.

The authors are with the Department of Electrical Engineering, Princeton University, Princeton, NJ 08540 USA (email: jpredd/kulkarni/poor@princeton.edu)

Consistency in Models for Distributed Learning under Communication Constraints¹

I. INTRODUCTION

A. Models for Distributed Learning

Consider the following learning model: Let X and Y be \mathcal{X} -valued and \mathcal{Y} -valued random variables, respectively, with a joint distribution denoted by \mathbf{P}_{XY} . \mathcal{X} is known as the feature, input, or observation space; \mathcal{Y} is known as the label, output, or target space. Throughout, we take $\mathcal{X} \subseteq \mathbb{R}^d$ and consider two cases corresponding to binary classification ($\mathcal{Y} = \{0, 1\}$) and regression estimation ($\mathcal{Y} = \mathbb{R}$). Given a loss function $l : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$, the decision-theoretic problem is to design a decision rule $g : \mathcal{X} \rightarrow \mathcal{Y}$ that achieves the minimal expected loss $L^* = \inf_g \mathbf{E}\{l(g(X), Y)\}$. Without prior knowledge of the distribution \mathbf{P}_{XY} , computing a loss minimizing decision rule is not possible. Instead, $D_n = \{(X_i, Y_i)\}_{i=1}^n$, an independent and identically distributed (i.i.d.) collection of training data with $(X_i, Y_i) \sim \mathbf{P}_{XY}$ for all $i \in \{1, \dots, n\}$ is available; the learning problem is to use this data to infer decision rules with small expected loss.

This standard learning model invites one to consider numerous questions; however in this work, we focus on the statistical property known as *universal consistency* [7], [12]. In traditional, centralized settings, D_n is provided to a single learning agent, and questions have been posed about the existence of classifiers or estimators that are universally consistent. The answers to such questions are well understood and are provided by results such as Stone's Theorem [35], [7], [12] and numerous others in the literature.

Suppose, in contrast with the standard centralized setting, that for each $i \in \{1, \dots, n\}$, the training datum (X_i, Y_i) is received by a distinct member of a network of n simple learning agents. When a central authority observes a new observation $X \sim \mathbf{P}_X$, it broadcasts the observation to the network in a request for information. At this time, each agent can respond with at most one bit. That is, each learning agent chooses whether or not to respond to the central authority's request for information; if it chooses to respond, an agent sends either a 1 or a 0 based on its local decision algorithm. Upon observing the

¹This paper was presented in part at the 17th Annual Conference on Learning Theory (COLT), Banff, Canada, July 1-4, 2004 [31] and in part at the 42nd Annual Allerton Conference on Communication, Control, and Computing, Monticello, IL, Sept 29-Oct 1, 2004 [32].

response of the network, the central authority acts as a fusion center, combining the information to create an estimate of Y . As in the centralized setting, a key question arises: do there exist agent decision rules and a fusion rule that result in a universally consistent network in the limit as the number of agents increases without bound?

In what follows, we answer this question in the affirmative for both binary classification and regression estimation. In the binary classification setting, we demonstrate agent decision rules and a fusion rule that correspond nicely with classical kernel classifiers. With this connection to classical work, the universal Bayes-risk consistency of this ensemble then follows immediately from celebrated analyses like Stone's Theorem, etc. In the regression setting, we demonstrate that under regularity, randomized agent decision rules exist such that when the central authority applies a scaled average vote combination of the agents' responses, the resulting estimator is universally consistent under L_2 -loss.

In this model, the agents convey slightly more information than is suggested by the mere one bit that we have allowed them to physically transmit to the fusion center. Indeed, each agent decides not between sending 1 or 0. Rather, each agent's decision rule can be viewed as a selection of one of *three* states: abstain, vote and send 0, and vote and send 1. With this observation, these results can be interpreted as follows: $\log_2(3)$ bits per agent per classification is sufficient for universal consistency to hold for both distributed classification and regression *with abstention*.

In this view, it is natural to ask whether these $\log_2(3)$ bits are necessary. Can consistency results be proven at lower bit rates? Consider a revised model, precisely the same as above, except that in response to the central authority's request for information, each agent must respond with 1 or 0; abstention is not an option and thus, each agent responds with exactly one bit per classification. Are there rules for which universal consistency results hold in distributed classification and regression *without abstention*?

Interestingly, we demonstrate that in the binary classification setting, randomized agent decision rules exist such that when a majority vote fusion rule is applied, universal Bayes-risk consistency holds. Next, we establish natural regularity conditions for candidate fusion rules and specify a reasonable class of agent decision rules. As an important negative result, we then demonstrate that for any agent decision rule within the class, there does not exist a regular fusion rule that is L_2 consistent for every distribution \mathbf{P}_{XY} . This result establishes the impossibility of universal consistency in this model for distributed regression without abstention for a restricted, but reasonable class of decision rules.

B. Motivation and Background

Motivation for studying distributed learning in general and the current models in particular arise from wireless sensor networks and distributed databases, applications that have attracted considerable attention in recent years [1]. Research in wireless sensor networks has focused on two separate aspects: networking issues, such as capacity, delay, and routing strategies; and applications issues. This paper is concerned with the second of these aspects, and in particular with the problem of distributed inference. Wireless sensor networks are *a fortiori* designed for the purpose of making inferences about the environments that they are sensing, and they are typically characterized by limited communications capabilities due to tight energy and bandwidth limitations, as well as the typically ad-hoc nature of wireless networks. Thus, distributed inference is a major issue in the study of wireless sensor networks.

In problems of distributed databases, there is a collection of training data that is massive in both the dimension of the feature space and quantity of data. For political, economic, social or technological reasons, this database is distributed geographically or in such a way that it is infeasible for any single agent to access the entire database. Multiple agents may be deployed to make inferences from various segments of the database, but communication constraints arising from privacy or security concerns highlight distributed inference as a key issue in this setting as well. Recent research has studied inference in the distributed databases setting from an algorithmic point of view; for example, [22] proposed a distributed boosting algorithm and studied its performance empirically.

Distributed detection and estimation is a well-developed field with a rich history. Much of the work in this area has focused on either parametric problems, in which strong statistical assumptions are made [36], [37], [3], [38], [23], [21], [6], [17], [8], or on traditional nonparametric formalisms, such as constant-false-alarm-rate detection [2]. Recently, [34] advocated a learning theoretic approach to wireless sensor networks and [26], in the context of kernel methods commonly used in machine learning, considered the classical model for decentralized detection [36] in a nonparametric setting.

In this paper, we consider an alternative nonparametric approach to the study of distributed inference that is most closely aligned with models considered in nonparametric statistics and the study of kernel estimators and other Stone-type rules. Extensive work has been done related to the consistency of Stone-type rules under various sampling processes; for example, [7], [12] and references therein, [5], [11], [18], [19], [20], [25], [27], [28], [29], [33], [35], [39], [40]. These models focus on various dependency structures within the training data and assume that a single processor has access to the entire data stream.

The nature of the work considered in this paper is to consider similar questions of universal consistency

in models that capture some of the structure in a distributed environment. As motivated earlier, agents in distributed scenarios have constrained communication capabilities and moreover, each may have access to distinct data streams that differ in distribution and may depend on parameters such as the state of a sensor network or location of a database. We consider the question: for a given model of communication amongst agents, each of whom has been allocated a small portion of a larger learning problem, can enough information can be exchanged to allow for a universally consistent ensemble? In this work, the learning problem is divided amongst agents by allocating each a unique observation of an i.i.d. sampling process. As explained earlier, we consider simple communication models with and without abstention. Insofar as these models present a useful picture of distributed scenarios, this paper addresses the issue of whether or not the guarantees provided by Stone’s Theorem in centralized environments hold in distributed settings. Notably, the models under consideration will be similar in spirit to their classical counterparts; indeed, similar techniques can be applied to prove results.

Note that [30] studies a similar model for distributed learning under communication constraints. Whereas [30] allocates regions of feature space amongst agents, here we allocate observations of an i.i.d. sampling process. Moreover, here we study a richer class of communication constraints. A related area of research lies in the study of ensemble methods in machine learning; examples of these techniques include bagging, boosting, mixtures of experts, and others [13], [4], [9], [10], [15]. These techniques are similar to the problem of interest here in that they aggregate many individually trained classifiers. However, the focus of these works is on the statistical and algorithmic advantages of learning with an ensemble and not on the nature of learning under communication constraints. Notably, [14] considered an PAC-like model for learning with many individually trained hypotheses in a distribution-specific (i.e., parametric) framework.

Numerous other works in the literature are relevant to the research presented here. However, different points need to be made depending on whether we consider regression or classification with or without abstention. Lacking such context here, we will save such discussion of these results for the appropriate sections in the paper.

C. Organization

The remainder of this paper is organized as follows. In Section II, the notation and technical assumptions relevant to the remainder of the paper are introduced. In Sections III and IV, we study the models for binary classification in communication with and without abstention, respectively. In Sections V and VI, we study the models for regression estimation with and without abstention in turn. In each section, we

present the main results, discuss important connections to other work in nonparametric statistics, and then proceed with a proof that further emphasizes differences from classical analyses like Stone's Theorem. In Section VII, we conclude with a discussion of future work. Technical lemmas that are readily apparent from the literature are left to the appendix.

II. PRELIMINARIES

In this section, we introduce notation and technical assumptions relevant to the remainder of the paper.

As stated earlier, let X and Y be \mathcal{X} -valued and \mathcal{Y} -valued random variables, respectively, with a joint distribution denoted by \mathbf{P}_{XY} . \mathcal{X} is known as the feature, input, or observation space; \mathcal{Y} is known as the label, output, or target space. Throughout, we will take $\mathcal{X} \subseteq \mathbb{R}^d$ and consider two cases corresponding to binary classification ($\mathcal{Y} = \{0, 1\}$) and regression estimation ($\mathcal{Y} = \mathbb{R}$). Let $D_n = \{(X_i, Y_i)\}_{i=1}^n$ denote an i.i.d. collection of training data with $(X_i, Y_i) \sim \mathbf{P}_{XY}$ for all $i \in \{1, \dots, n\}$.

Throughout this paper, we will use δ_{ni} to denote the randomized response of the i^{th} learning agent in an ensemble of n agents. For each $i \in \{1, \dots, n\}$, δ_{ni} is an \mathcal{S} -valued random variable, where \mathcal{S} is the decision space for the agent; in models *with abstention* we take $\mathcal{S} = \{\text{abstain}, 1, 0\}$ and in models *without abstention* we take $\mathcal{S} = \{1, 0\}$. As an important consequence of the assumed lack of inter-agent communication and the assumption that D_n is i.i.d., we have the following observation which will be fundamental to the subsequent analysis:

- (A) The i^{th} agent's response, δ_{ni} , may be dependent on X , X_i , and Y_i , but is statistically independent of $\{(X_j, Y_j)\}_{j \neq i}$ and conditionally independent of $\{\delta_{nj}\}_{j \neq i}$ given X .

Thus, to specify δ_{ni} and thereby design agent decision rules, it suffices to define the conditional distribution $\mathbf{P}\{\delta_{ni} | X, X_i, Y_i\}$ for all $(X, X_i, Y_i) \in \mathcal{X} \times \mathcal{X} \times \mathcal{Y}$. In each of the subsequent sections, we will find it convenient to do so by specifying a function $\bar{\delta}_n(x) : \mathcal{X} \times \mathcal{X} \times \mathcal{Y} \rightarrow \{\text{abstain}\} \cup [0, 1]$. In particular, we define

$$\begin{aligned} \mathbf{P}\{\delta_{ni} = \text{abstain} | X, X_i, Y_i\} &= \begin{cases} 1, & \text{if } \bar{\delta}_n(X, X_i, Y_i) = \text{abstain} \\ 0, & \text{otherwise} \end{cases} \\ \mathbf{P}\{\delta_{ni} = 1 | X, X_i, Y_i\} &= \begin{cases} 0, & \text{if } \bar{\delta}_n(X, X_i, Y_i) = \text{abstain} \\ \bar{\delta}_n(X, X_i, Y_i), & \text{otherwise} \end{cases} \\ \mathbf{P}\{\delta_{ni} = 0 | X, X_i, Y_i\} &= \begin{cases} 0, & \text{if } \bar{\delta}_n(X, X_i, Y_i) = \text{abstain} \\ 1 - \bar{\delta}_n(X, X_i, Y_i), & \text{otherwise} \end{cases} \end{aligned} \quad (1)$$

It is straightforward to verify that (1) is a valid probability distribution for every $(X, X_i, Y_i) \in \mathcal{X} \times \mathcal{X} \times \mathcal{Y}$. Therefore, together with (A), δ_{ni} is clearly specified by $\bar{\delta}_{ni}(x)$ and (1).

Note, this formalism serves merely as a technical convenience and should not mask the simplicity of the agent decision rules. In words, an agent will abstain from voting if $\bar{\delta}_n(X, X_i, Y_i) = \text{abstain}$; else, the agent flips a biased coin to send 1 or 0, with the bias determined by $\bar{\delta}_n(X, X_i, Y_i)$. Though this formalism may appear restrictive since rules of this form do not allow randomized decisions to abstain, the results in this paper do not rely on this flexibility.

To emphasize, note that communication is constrained between the agents and the fusion center via the limited decision space \mathcal{S} and as above, communication between agents is not allowed (the latter is a necessary precondition for observation (A)). Consistent with the notation, we assume that the agents have knowledge of n , the number of agents in the ensemble. Moreover, we assume that for each n , every agent has the same local decision rule; i.e., the ensemble is homogenous in this sense. An underlying assumption is that each agent is able to generate random numbers, independent of the rest of the network.

Consistent with convention, we use $g_n(x) = g_n(x, \{\delta_{ni}\}_{i=1}^n) : \mathcal{X} \times \mathcal{S}^n \rightarrow \{0, 1\}$ to denote the central authority's fusion rule in the binary classification frameworks and similarly, we use $\hat{\eta}_n(x) = \hat{\eta}_n(x, \{\delta_{ni}\}_{i=1}^n) : \mathcal{X} \times \mathcal{S}^n \rightarrow \mathbb{R}$ to denote its fusion rule in the regression frameworks. In defining fusion rules throughout the remainder of the paper, it will be convenient to denote the random set $I_V = I_V(X, D_n) \triangleq \{i \in \{1, \dots, n\} : \delta_{ni} \neq \text{abstain}\}$ as the set of agents that vote and hence, do not abstain. To emphasize the central authority's primary role of aggregating the response of the network, we shall henceforth refer to this agent as a *fusion center*.

Defining a loss function $l : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$, we seek ensembles that achieve the minimal expected loss. In the binary classification setting, the criterion of interest is the probability of misclassification; we let $l(y, y') = 1_{\{y \neq y'\}}$, the well-known zero-one loss. The structure of the risk minimizing MAP decision rule is well-understood [7]; let $\delta_B : \mathcal{X} \rightarrow \{0, 1\}$ denote this Bayes decision rule. In regression settings, we consider the squared error criterion; we let $l(y, y') = |y - y'|^2$. It is well known that the regression function

$$\eta(x) = \mathbf{E}\{Y | X = x\} \tag{2}$$

achieves the minimal expected loss in this case. Throughout the remainder of the paper, we let $L^* = \inf_f \mathbf{E}\{l(f(X), Y)\}$ denote the minimal expected loss. Depending on whether we find ourselves in the binary classification or regression setting, it will be clear from the context whether L^* refers to the optimal (binary) Bayes risk or minimal mean squared error.

In this work, we focus on the statistical property known as *universal consistency* [7], [12], defined as follows.

Definition 1: Let $L_n = \mathbf{E}\{l(f_n(X, D_n), Y) | D_n\}$. $\{f_n\}_{n=1}^{\infty}$ is said to be *universally consistent* if $\mathbf{E}\{L_n\} \rightarrow L^*$ for all distributions \mathbf{P}_{XY} .

This definition requires convergence in expectation and according to convention, defines *weak* universal consistency. This notion is contrasted with *strong* universal consistency where $L_n \rightarrow L^*$ almost surely. Extending results of weak universal consistency to the strong sense has generally required the theory of large deviations, in particular McDiarmid's inequality [7]. Though the focus in this paper is on the weaker sense, the results in this paper might be extended to strong universal consistency using similar techniques. In particular, note that since consistency in distributed classification *with abstention* can be reduced to Stone's Theorem, the extension to strong universal consistency follows immediately from standard results. Further, the negative result for distributed regression *without abstention* automatically precludes consistency in the strong sense. An extension for distributed classification without abstention and distributed regression with abstention may be possible under a refined analysis; the authors leave such analysis for future research.

III. DISTRIBUTED CLASSIFICATION WITH ABSTENTION: STONE'S THEOREM

In this section, we show that the universal consistency of distributed classification with abstention follows immediately from Stone's Theorem and the classical analysis of naive kernel classifiers. To start, let us briefly recap the model. Since we are in the classification framework, $\mathcal{Y} = \{0, 1\}$. Suppose that for each $i \in \{1, \dots, n\}$, the training datum $(X_i, Y_i) \in D_n$ is received by a distinct member of a network of n learning agents. When the fusion center observes a new observation $X \sim \mathbf{P}_X$, it broadcasts the observation to the network in a request for information. At this time, each of the learning agents can respond with at most one bit. That is, each learning agent chooses whether or not to respond to the fusion center's request for information; and if an agent chooses to respond, it sends either a 1 or a 0 based on a local decision algorithm. Upon receiving the agents' responses, the fusion center combines the information to create an estimate of Y .

To answer the question of whether agent decision rules and fusion rules exist that result in a universally consistent ensemble, let us construct one natural choice. With $B_{r_n}(x) = \{x' \in \mathbb{R}^d : \|x - x'\|_2 \leq r_n\}$, let

$$\bar{\delta}_n(x, X_i, Y_i) = \begin{cases} Y_i, & \text{if } X_i \in B_{r_n}(x) \\ \text{abstain,} & \text{otherwise} \end{cases} \quad (3)$$

and

$$g_n(x) = \begin{cases} 1, & \text{if } \sum_{i \in I_V} \delta_{ni} \geq \frac{1}{2} |I_V| \\ 0, & \text{otherwise} \end{cases}, \quad (4)$$

so that $g_n(x)$ amounts to a majority vote fusion rule. Recall from (1) that the agents' randomized responses are defined by $\bar{\delta}_n(\cdot)$. In words, agents respond according to their training data label as long as the new observation X is sufficiently close to their training observation X_i ; else, they abstain. In this model with abstention, note that δ_{ni} is $\{\text{abstain}, 1, 0\}$ -valued since Y_i is binary valued and thus, the communications constraints are obeyed.

With this choice, it is straightforward to see that the net decision rule is equivalent to the plug-in kernel classifier rule with the naive kernel. Indeed,

$$g_n(x) = \begin{cases} 1, & \text{if } \frac{\sum_{i=1}^n Y_i 1_{B_{r_n}(x)}(X_i)}{\sum_{i=1}^n 1_{B_{r_n}(x)}(X_i)} \geq \frac{1}{2} \\ 0, & \text{otherwise} \end{cases}. \quad (5)$$

With this equivalence², the universal consistency of the ensemble follows from Stone's Theorem applied to naive kernel classifiers. With $L_n = \mathbf{P}\{g_n(X) \neq Y | D_n\}$, the probability of error of the ensemble conditioned on the random training data, we state this known result without proof as Theorem 1.

Theorem 1: ([7]) If $r_n \rightarrow 0$ and $(r_n)^d n \rightarrow \infty$ as $n \rightarrow \infty$, then $\mathbf{E}\{L_n\} \rightarrow L^*$ for all distributions \mathbf{P}_{XY} .

The kernel classifier with the naive kernel is somewhat unique amongst other frequently analyzed universally consistent classifiers in its relevance to the current model. More general kernels (for instance, a Gaussian kernel) are not easily applicable as the real-valued weights do not naturally form a randomized decision rule. Furthermore, nearest neighbor rules do not apply as a given agent's decision rule would then need to depend on the data observed by the other agents; such inter-agent communication is not allowed in the current model.

IV. DISTRIBUTED CLASSIFICATION WITHOUT ABSTENTION

As noted in the introduction, given the result of the previous section, it is natural to ask whether the communication constraints can be tightened. Let us consider the second model in which the agents cannot choose to abstain. In effect, each agent communicates one bit per decision. Again, we consider the binary classification framework but as a technical convenience, adjust our notation so that $\mathcal{Y} = \{+1, -1\}$ instead of the usual $\{0, 1\}$; also, agents now decide between sending ± 1 . The formalism introduced in

²Strictly speaking, this equality holds almost surely (a.s.), since the agents' responses are random variables.

Section II can be extended naturally to allow this slight modification; we allow δ_{ni} to be specified so that $\mathbf{P}\{\delta_{ni} = +1 | X, X_i, Y_i\} = \bar{\delta}_{ni}(x, X_i, Y_i)$. We again consider whether universally Bayes-risk consistent schemes exist for the ensemble.

Consider the randomized agent decision rule specified as follows:

$$\bar{\delta}_{ni}(x, X_i, Y_i) = \begin{cases} \frac{1}{2}Y_i + \frac{1}{2}, & \text{if } X_i \in B_{r_n}(x) \\ \frac{1}{2}, & \text{otherwise} \end{cases}. \quad (6)$$

Recall from (1) that the agents' randomized responses are defined by $\bar{\delta}_n(\cdot)$. Note that $\mathbf{P}\{\delta_{ni} = Y_i | X_i \in B_{r_n}(x)\} = 1$, and thus, the agents respond according to their training data label if x is sufficiently close to X_i . Else, they simply “guess”, flipping an unbiased coin. In this model without abstention, it is readily verified that each agent transmits one bit per decision as δ_{ni} is $\{\pm 1\}$ -valued since $\mathbf{P}\{\delta_{ni} = \text{abstain}\} = 0$; thus, the communication constraints are obeyed.

A natural fusion rule is the majority vote. That is, the fusion center decides according to

$$g_n(x) = \begin{cases} 1, & \text{if } \sum_{i=1}^n \delta_{ni} > 0 \\ -1, & \text{otherwise} \end{cases}. \quad (7)$$

As before, the natural performance metric for the ensemble is the probability of misclassification. Modifying our convention slightly, let $D_n = \{(X_i, Y_i, \delta_{ni})\}_{i=1}^n$ and define

$$L_n = \mathbf{P}\{g_n(X) \neq Y | D_n\}. \quad (8)$$

That is, L_n is the conditional probability of error of the majority vote fusion rule conditioned on the randomness in agent training and agent decision rules.

A. Main Result and Comments

Theorem 2 specifies sufficient conditions for consistency for an ensemble using the described decision rules.

Theorem 2: If $r_n \rightarrow 0$ and $(r_n)^d \sqrt{n} \rightarrow \infty$ as $n \rightarrow \infty$, then $\mathbf{E}\{L_n\} \rightarrow L^*$.

Yet again, the conditions of the theorem strike a similarity with consistency results for kernel classifiers using the naive kernel. Indeed, $r_n \rightarrow 0$ ensures that the bias of the classifier decays to zero. However, $\{r_n\}_{n=1}^\infty$ must not decay too rapidly. As the number of agents in the ensemble grows large, many, indeed most, of the agents will be “guessing” for any given classification; in general, only a decaying fraction of the agents will respond with useful information. In order to ensure that these informative bits can be heard through the noise introduced by the guessing agents, $(r_n)^d \sqrt{n} \rightarrow \infty$. Note the difference between

this result and that for naive kernel classifiers where $(r_n)^d n \rightarrow \infty$ assures a sufficient rate of convergence for $\{r_n\}_{n=1}^\infty$.

Notably, to prove this result, we show directly that the expected probability of misclassification converges to the Bayes rate. This is unlike techniques commonly used to demonstrate the consistency of kernel classifiers, etc., which are so-called “plug-in” classification rules. These rules estimate the *a posteriori* probabilities $\mathbf{P}\{Y = i | X\}$, $i = \pm 1$ and construct classifiers based on thresholding the estimate. In this setting, it suffices to show that these estimates converge to the true probabilities in $L^p(\mathbf{P}_X)$. However, for this model, we cannot estimate the *a posteriori* probabilities and must resort to another proof technique; this foreshadows the negative result of Section VI.

With our choice of “coin flipping” agent decision rules, one may be tempted to model the observations made by the fusion center as noise-corrupted labels from the training set and to thereby recover Theorem 2 from the literature on learning with noisy data. However, note that since the fusion center does not have access to the agents’ feature observations (i.e., $\{X_i\}_{i=1}^n$), the fusion rule cannot in general be modeled as a “plug-in” classification rule as analyzed, for instance, in [24]. Moreover, in contrast to the noise models considered in [24], the agent decision rules here are statistically dependent on X and are also dependent on X_i in an atypical way: the noise statistics depend on n and for particular \mathbf{P}_{XY} , one can show that as n increases without bound, the probability that an agent guesses (a label is noisy) grows toward 1. These differences distinguish Theorem 2 from results in the literature on learning with noisy data.

B. Proof of Theorem 2

Proof: Fix an arbitrary $\epsilon > 0$. We will show that $\mathbf{E}\{L_n\} - L^*$ is less than ϵ for all sufficiently large n . Using the notation in (2), we write $\eta(x) = \mathbf{E}\{Y | X = x\} = \mathbf{P}\{Y = +1 | X = x\} - \mathbf{P}\{Y = -1 | X = x\}$ and define $A_\epsilon = \{x : |\eta(x)| > \frac{\epsilon}{2}\}$. It follows that

$$\begin{aligned} \mathbf{E}\{L_n\} - L^* &= \mathbf{E}\left\{\mathbf{P}\{g_n(X) \neq Y | D_n\}\right\} - \mathbf{P}\{\delta_B(X) \neq Y\} \\ &= \mathbf{E}\left\{\left(\mathbf{P}\{g_n(X) \neq Y | D_n, X\} - \mathbf{P}\{\delta_B(X) \neq Y | X\}\right)\left(1_{A_\epsilon}(X) + 1_{\bar{A}_\epsilon}(X)\right)\right\}, \quad (9) \end{aligned}$$

with the expectation in (9) being taken with respect to X and D_n . Note that for all $x \in \bar{A}_\epsilon$, $\mathbf{P}\{\delta_B(X) \neq Y | X = x\} = \frac{1}{2} - \frac{|\eta(x)|}{2} \geq \frac{1}{2} - \frac{\epsilon}{4}$ and therefore, $\mathbf{P}\{g_n(X) \neq Y | D_n, X\} \leq 1 - \mathbf{P}\{\delta_B(X) \neq Y | X = x\} \leq \frac{1}{2} + \frac{\epsilon}{4}$. Thus,

$$\begin{aligned} \mathbf{E}\{L_n\} - L^* &\leq \mathbf{E}\left\{\left(\mathbf{P}\{g_n(X) \neq Y | D_n, X\} - \mathbf{P}\{\delta_B(X) \neq Y | X\}\right)1_{A_\epsilon}(X) + \frac{\epsilon}{2}\right\} \\ &\leq \mathbf{P}\left\{g_n(X) \neq \delta_B(X) \mid X \in A_\epsilon\right\}\mathbf{P}\left\{A_\epsilon\right\} + \frac{\epsilon}{2}. \end{aligned}$$

Note that if $\mathbf{P}\{A_\epsilon\} = 0$, then the proof is complete. Let us proceed assuming $\mathbf{P}\{A_\epsilon\} > 0$. Clearly, it suffices to show that $\lim_{n \rightarrow \infty} \mathbf{P}\left\{g_n(X) \neq \delta_B(X) \mid X \in A_\epsilon\right\} \leq \frac{\epsilon}{2}$. Let us define the quantities

$$m_n(x) = \mathbf{E}\{\eta(X)\delta_{ni} \mid X = x\}$$

$$\sigma_n^2(x) = \mathbf{E}\{|\eta(X)\delta_{ni} - m_n(X)|^2 \mid X = x\},$$

with the expectation being taken over the random training data and the randomness introduced by the agent decision rules. Respectively, $m_n(x)$ and $\sigma_n^2(x)$ can be interpreted as the mean and variance of the “margin” of the agent response δ_{ni} , conditioned on the observation X . For large positive $m_n(x)$, the agents can be expected to respond “confidently” (with large margin) according to the Bayes rule when asked to classify an object x . For large $\sigma_n^2(x)$, the fusion center can expect to observe a large variance amongst the individual agent responses to x .

Fix any integer $k > 0$. Consider the sequence of sets indexed by n ,

$$B_{n,k} = \{x \in \mathcal{X} : m_n(x)n > k\sqrt{n}\sigma_n(x)\},$$

so that $x \in B_{n,k}$ if and only if $\frac{m_n(x)\sqrt{n}}{\sigma_n(x)} > k$. We can interpret $B_{n,k}$ as the set of observations for which informed agents have a sufficiently strong signal compared with the noise of the guessing agents. Then,

$$\begin{aligned} \mathbf{P}\left\{g_n(X) \neq \delta_B(X) \mid X \in A_\epsilon\right\} &= \mathbf{P}\left\{\eta(X) \sum_{i=1}^n \delta_{ni} < 0 \mid X \in A_\epsilon\right\} \\ &= \mathbf{P}\left\{\eta(X) \sum_{i=1}^n \delta_{ni} < 0 \mid X \in A_\epsilon \cap B_{n,k}\right\} \mathbf{P}\{X \in B_{n,k} \mid X \in A_\epsilon\} + \\ &\quad \mathbf{P}\left\{\eta(X) \sum_{i=1}^n \delta_{ni} < 0 \mid X \in A_\epsilon \cap \bar{B}_{n,k}\right\} \mathbf{P}\{X \in \bar{B}_{n,k} \mid X \in A_\epsilon\} \end{aligned} \quad (10)$$

Note that conditioned on X , $\eta(X) \sum_{i=1}^n \delta_{ni}$ is a sum of independent and identically distributed random variables with mean $m_n(X)$ and variance $\sigma_n^2(X)$. Further, for $x \in B_{n,k}$, $\eta(x) \sum_{i=1}^n \delta_{ni} < 0$ implies $|\eta(x) \sum_{i=1}^n \delta_{ni} - m_n(x)n| > k\sqrt{n}\sigma_n^2(x)$. Thus, it is straightforward to see that,

$$\begin{aligned} \mathbf{P}\left\{\eta(X) \sum_{i=1}^n \delta_{ni} < 0 \mid X \in A_\epsilon \cap B_{n,k}\right\} &= \mathbf{E}\left\{\mathbf{P}\left\{\eta(X) \sum_{i=1}^n \delta_{ni} < 0 \mid X\right\} \mid X \in A_\epsilon \cap B_{n,k}\right\} \\ &\leq \mathbf{E}\left\{\mathbf{P}\left\{\left|\eta(X) \sum_{i=1}^n \delta_{ni} - m_n(X)n\right| > k\sqrt{n}\sigma_n(X) \mid X\right\} \mid X \in A_\epsilon \cap B_{n,k}\right\} \\ &\leq \frac{1}{k^2}. \end{aligned}$$

Here, the last statement follows from Markov's Inequality. Choosing k sufficiently large and returning to (10),

$$\mathbf{P}\left\{g_n(X) \neq \delta_B(X) \mid X \in A_\epsilon\right\} \leq \frac{\epsilon}{2} + \mathbf{P}\{X \in \bar{B}_{n,k} \mid X \in A_\epsilon\}.$$

Now let us determine specific expressions for $m_n(x)$ and $\sigma_n^2(x)$, as dictated by our choice of agent decision rules. Clearly,

$$\begin{aligned} m_n(x) &= \eta(x)\mathbf{E}\{\delta_{ni} \mid X = x\} \\ &= \eta(x)\mathbf{E}\left\{\mathbf{E}\{2\bar{\delta}_{ni}(X, X_i, Y_i) - 1 \mid X, X_i, Y_i\} \mid X = x\right\} \\ &= \eta(x)\left(0 \cdot \mathbf{P}\{X_i \in \bar{B}_{r_n}(x)\} + \eta_n(x) \cdot \mathbf{P}\{X_i \in B_{r_n}(x)\}\right) \\ &= \eta(x)\eta_n(x) \int 1_{B_{r_n}(x)}(y)P_X(dy), \end{aligned}$$

with $\eta_n(x) = \mathbf{E}\{\eta(X) \mid X \in B_{r_n}(x)\}$. Also,

$$\begin{aligned} \sigma_n^2(x) &= \eta^2(x)\mathbf{E}\{|\delta_{ni} - \mathbf{E}\{\delta_{ni} \mid X = x\}|^2 \mid X = x\} \\ &= \eta^2(x)(1 - \mathbf{E}\{\delta_{ni} \mid X = x\}^2). \end{aligned}$$

Thus,

$$\begin{aligned} \mathbf{P}\{X \in \bar{B}_{n,k} \mid X \in A_\epsilon\} &= \mathbf{P}\{m_n(X)n < k\sqrt{n}\sigma_n(X) \mid X \in A_\epsilon\} \\ &= \mathbf{P}\left\{\frac{\eta(X)\eta_n(X) \int 1_{B_{r_n}(X)}(y)P_X(dy)\sqrt{n}}{|\eta(X)|\sqrt{1 - \mathbf{E}\{\delta_{ni} \mid X\}^2}} < k \mid X \in A_\epsilon\right\} \\ &= \mathbf{P}\left\{\left(\text{sgn}(\eta(X))\eta_n(X)\right)\left(\frac{\sqrt{n}}{\sqrt{1 - \mathbf{E}\{\delta_{ni} \mid X\}^2}} \int 1_{B_{r_n}(X)}(y)P_X(dy)\right) < k \mid X \in A_\epsilon\right\}. \end{aligned}$$

For any $1 \geq \gamma > 0$, we have

$$\begin{aligned} \mathbf{P}\{X \in \bar{B}_{n,k} \mid X \in A_\epsilon\} &\leq \mathbf{P}\left\{\frac{\sqrt{n}}{\sqrt{1 - \mathbf{E}\{\delta_{ni} \mid X\}^2}} \int 1_{B_{r_n}(X)}(y)P_X(dy) < k \mid \right. \\ &\quad \left. X \in A_\epsilon, \text{sgn}(\eta(X))\eta_n(X) > \gamma\right\} + \mathbf{P}\{\text{sgn}(\eta(X))\eta_n(X) \leq \gamma \mid X \in A_\epsilon\}. \end{aligned} \quad (11)$$

First, consider the second term. With $\gamma = \frac{\epsilon}{4}$, it follows from our choice of A_ϵ that $\{\text{sgn}(\eta(X))\eta_n(X) \leq \frac{\epsilon}{4}\}$ implies $\{|\eta(X) - \eta_n(X)| > \frac{\epsilon}{4}\}$. Thus,

$$\mathbf{P}\left\{\text{sgn}(\eta(X))\eta_n(X) \leq \frac{\epsilon}{4} \mid X \in A_\epsilon\right\} \leq \mathbf{P}\left\{|\eta(X) - \eta_n(X)| > \frac{\epsilon}{4} \mid X \in A_\epsilon\right\}.$$

Since by technical Lemma 2 (see appendix), $\eta_n(X) \rightarrow \eta(X)$ in probability and by assumption $\mathbf{P}\{A_\epsilon\} > 0$, it follows from technical Lemma 1 in the appendix that $\mathbf{P}\{\text{sgn}(\eta(X))\eta_n(X) \leq \frac{\epsilon}{4} \mid X \in A_\epsilon\} \rightarrow 0$.

Returning to (11) with $\gamma = \frac{\epsilon}{4}$, note that we have just demonstrated that $\lim_{n \rightarrow \infty} \mathbf{P}\{\text{sgn}(\eta(X))\eta_n(X) > \frac{\epsilon}{4}\} = 1$. Thus, to show that the first term converges to zero, by technical Lemma 1, it suffices to show that

$$\frac{\sqrt{n}}{\sqrt{1 - \mathbf{E}\{\delta_{ni} | X\}^2}} \int 1_{B_{r_n}(X)}(y) P_X(dy) \rightarrow \infty \text{ i.p.} \quad (12)$$

Since $\frac{1}{\sqrt{1 - \mathbf{E}\{\delta_{ni} | X\}^2}} \geq 1$, this follows from technical Lemma 3 in the appendix and the fact that $(r_n)^d \sqrt{n} \rightarrow \infty$. This completes the proof. \blacksquare

V. DISTRIBUTED REGRESSION WITH ABSTENTION

We now turn our attention to distributed regression. As in Section III, the model remains the same except that now $\mathcal{Y} = \mathbb{R}$; that is, Y is an \mathbb{R} -valued random variable and likewise, agents receive real-valued training data labels, Y_i . In this section, we consider communication with abstention. With the aim of determining whether universally consistent ensembles can be constructed, let us devise candidate rules.

For some as yet unspecified sequence of functions $T_n : \mathbb{R} \rightarrow [0, 1]$ and a sequence of real numbers $\{r_n\}_{n=1}^\infty$, consider the randomized agent decision rules specified as follows:

$$\bar{\delta}_{ni}(x) = \begin{cases} T_n(Y_i) & \text{if } X_i \in B_{r_n}(x) \\ \text{abstain,} & \text{otherwise} \end{cases}, \quad (13)$$

for $i = 1, \dots, n$. In words, the agents choose to vote only if X_i is close enough to X ; to vote, they flip a biased coin, with the bias determined by the size of the ensemble n and Y_i , via the function $T_n(\cdot)$. In this model with abstention, note that δ_{ni} is $\{\text{abstain}, 1, 0\}$ -valued and thus, the communication constraints are obeyed.

It is intuitively clear that $T_n(\cdot)$ should be designed so that the realization of random bit $\delta_{n,i}$ reveals information about the real-valued label Y_i to the fusion center. In particular, it is natural to ask whether any continuous bijective mapping \mathbb{R} to the interval $(0, 1)$ would suffice in biasing the coin in a manner that is informative enough to provide universal consistency. For example, one might chose $T_n(y) = T(y) = \frac{1}{1+e^{-y}}$ and consider agent decision rules of the form (13) in conjunction with a fusion rule like

$$\hat{\eta}_n(x) = T^{-1}\left(\frac{\sum_{i \in I_V} \delta_{ni}}{|I_V|}\right). \quad (14)$$

Since agents have the flexibility to abstain, the fusion center can accurately estimate the average bias chosen by non-abstaining agents; the hope, then, is to determine the corresponding average label by

inverting $T(\cdot)$. As observed in the proof, such a choice is not possible, in general, since $T(\cdot)$ is nonlinear; such an approach introduces a systematic bias to the estimator and thereby prevents consistency.

If, however, $|Y| \leq B$ a.s. for some known $B > 0$, it suffices to choose $T_n(\cdot)$ as the linear function mapping $[-B, B]$ to $[0, 1]$. Since in this case, $T_n^{-1}(\mathbf{E}\{\delta_{n,i} | X, X_i\}) = \mathbf{E}\{Y_i | X_i\}$, universal consistency then follows with trivial modifications to the proof of Stone's Theorem.

This intuition leads us to a rule that captures consistency in the general case. Though choices abound, we can choose T_n to be piecewise linear. In particular, let $\{c_n\}_{n=1}^{\infty}$ be an arbitrary sequence of real numbers such that $c_n \rightarrow \infty$ as $n \rightarrow \infty$ and choose,

$$T_n(Y_i) = \begin{cases} \frac{1}{2c_n} Y_i + \frac{1}{2} & |Y_i| \leq c_n \\ \frac{1}{2}, & \text{otherwise} \end{cases}, \quad (15)$$

and specify the fusion rule as

$$\hat{\eta}_n(x) = 2c_n \left(\frac{\sum_{i \in I_V} \delta_{ni}}{|I_V|} - \frac{1}{2} \right). \quad (16)$$

In words, the fusion center shifts and scales the average vote. For appropriately chosen sequences $\{c_n\}_{n=1}^{\infty}$ and $\{r_n\}_{n=1}^{\infty}$, this ensemble is universally consistent, as proved by Theorem 3.

In particular, we will consider $L_n = \mathbf{E}\{|\hat{\eta}_n(X) - Y|^2\}$ with the expectation being taken over X , $D_n = \{(X_i, Y_i)\}_{i=1}^n$, and the randomness introduced in the agent decision rules.

A. Main Result and Comments

Assuming an ensemble using the described decision rules, Theorem 3 specifies sufficient conditions for consistency.

Theorem 3: Suppose \mathbf{P}_{XY} is such that \mathbf{P}_X is compactly supported and $\mathbf{E}\{Y^2\} < \infty$. If, as $n \rightarrow \infty$,

- 1) $c_n \rightarrow \infty$,
- 2) $r_n \rightarrow 0$, and
- 3) $\frac{c_n^2}{nr_n^d} \rightarrow 0$,

then $\mathbf{E}\{L_n\} \rightarrow L^*$.

More generally, the constraint regarding the compactness of \mathbf{P}_X can be weakened. As will be observed in the proof below, \mathbf{P}_X must be such that when coupled with a bounded random variable Y , there is a known convergence rate of the variance term of the naive kernel classifier (under a standard i.i.d. sampling model). $\{c_n\}_{n=1}^{\infty}$ should be chosen so that it grows at a rate slower than the rate at which the variance term decays. Notably, to select $\{c_n\}_{n=1}^{\infty}$, one does not need to understand the convergence rate

of the bias term, and this is why continuity conditions are not required; the bias term will converge to zero universally as long as $c_n \rightarrow \infty$ and $r_n \rightarrow 0$ as $n \rightarrow \infty$.

In observing the response of the network, the fusion center sees δ_{ni} from those agents who have not abstained. Since these random variables can be viewed as random quantizations or transformations of the labels in the training data, it is natural to ask whether the consistency of these rules follows as a special case of models for learning with noisy data. In this case, the underlying noise model would transform the label Y_i to the set $\{0, 1\}$ in a manner that would be statistically dependent on X , X_i , Y_i itself and n . Though it is possible to view the current question in this framework, to our knowledge such a highly structured noise model has not been considered in the literature.

Finally, those familiar with the classical statistical pattern recognition literature will find the style of proof very familiar; special care must be taken to demonstrate that the variance of the estimate does not decrease too slowly compared to $\{c_n\}_{n=1}^{\infty}$ and to show that the bias introduced by the ‘‘clipped’’ agent decision rules converges to zero.

B. Proof of Theorem 3

Proof: By standard orthogonality arguments [12], it suffices to show that $\mathbf{E}\{|\hat{\eta}_n(X) - \eta(X)|^2\} \rightarrow 0$ as $n \rightarrow 0$.

Define $\bar{\eta}_n(x) \triangleq \mathbf{E}\{\delta_{ni} | X_i = x, \|X - X_i\| \leq r_n\}$. Proceeding in the traditional manner, note that by the standard inequality

$$(a_1 + \dots + a_k)^2 \leq k(a_1^2 + \dots + a_k^2), \quad (17)$$

it follows that

$$\begin{aligned} \mathbf{E}\{|\hat{\eta}_n(X) - \eta(X)|^2\} &\leq 2\mathbf{E}\left\{\left|2c_n\left(\frac{\sum_{i \in I_V} \delta_{ni}}{|I_V|} - \frac{1}{2}\right) - 2c_n\left(\frac{\sum_{i \in I_V} \bar{\eta}_n(X_i)}{|I_V|} - \frac{1}{2}\right)\right|^2\right\} \\ &\quad + 2\mathbf{E}\left\{\left|2c_n\left(\frac{\sum_{i \in I_V} \bar{\eta}_n(X_i)}{|I_V|} - \frac{1}{2}\right) - \eta(X)\right|^2\right\} \\ &\triangleq J_n + K_n. \end{aligned}$$

Starting with the first term,

$$\begin{aligned} J_n &= 8c_n^2 \mathbf{E}\left\{\left|\frac{\sum_{i \in I_V} (\delta_{ni} - \bar{\eta}_n(X_i))}{|I_V|}\right|^2\right\} \\ &= 8c_n^2 \mathbf{E}\left\{\mathbf{E}\left\{\frac{\sum_{i \in I_V} (\delta_{ni} - \bar{\eta}_n(X_i))^2}{|I_V|^2} \middle| X, X_1, \dots, X_n\right\}\right\}. \end{aligned}$$

Here, the first equality follows from algebra; the second follows after noting that for all $i \in I_V$, $\mathbf{E}\{\delta_{ni} | X, X_1, \dots, X_n\} = \hat{\eta}_n(X_i)$ and canceling out cross-terms in the expansion of the squared sum in

the numerator. Note that conditioned on X and X_i , δ_{ni} is Bernoulli with parameter $\bar{\eta}_n(X_i)$ for all $i \in I_V$. Thus, bounding the variance of a Bernoulli random variable, we continue above,

$$\leq 2c_n^2 \mathbf{E} \left\{ \frac{1}{|I_V|} 1_{\{|I_V|>0\}} \right\}.$$

Here we have applied the convention $\frac{0}{0} = 0$. Conditioning on X and applying technical Lemma 4 (see the appendix) to the binomial random variable $|I_V| = \sum_{i=1}^n 1_{\{X_i \in B_{r_n}(X)\}}$, it follows that,

$$J_n \leq 2c_n^2 \mathbf{E} \left\{ \frac{2}{n \mathbf{P}_{X_1} \{X_1 \in B_{r_n}(X)\}} \right\}. \quad (18)$$

Here, for convenience, we have exploited the fact that D_n is i.i.d. and reused the variable X_1 . Since \mathbf{P}_X is compactly supported, the expectation in (18) can be bounded by a term $O(\frac{1}{nr_n^d})$ using an argument typically used to demonstrate the consistency of kernel estimators [12]. For completeness, we include it here.

Since S , the support of \mathbf{P}_X , is compact, we can find $z_1, \dots, z_{M_n} \in \mathbb{R}^d$ such that $S \subseteq \cup_{i=1}^{M_n} B_{r_n/2}(z_i)$ and $M_n \leq \frac{c_1}{r_n^d}$ for some constant c_1 . Thus,

$$\begin{aligned} 2c_n^2 \mathbf{E} \left\{ \frac{2}{n \mathbf{P}_{X_1} \{X_1 \in B_{r_n}(X)\}} \right\} &\leq 4c_n^2 \sum_{i=1}^{M_n} \mathbf{E} \left\{ \frac{1_{\{B_{r_n/2}(z_i)\}}(X)}{n \mathbf{P}_{X_1} \{X_1 \in B_{r_n}(X)\}} \right\} \\ &\leq 4c_n^2 \sum_{i=1}^{M_n} \mathbf{E} \left\{ \frac{1_{\{B_{r_n/2}(z_i)\}}(X)}{n \mathbf{P}_{X_1} \{X_1 \in B_{r_n/2}(z_i)\}} \right\} \\ &= \frac{4c_n^2 M_n}{n} \\ &\leq \frac{4c_1 c_n^2}{nr_n^d}. \end{aligned}$$

Finally, by condition (3) of Theorem 3, it follows that $J_n \rightarrow 0$. Note that J_n is essentially the variance of the estimator. Much of the work thus far has been the same as showing that in traditional i.i.d. sampling process settings, the variance of the naive kernel is universally bounded by a term $O(\frac{1}{nr_n^d})$ when \mathbf{P}_X is compactly supported and Y is bounded [12]. This observation is consistent with the comments above.

Now, let us consider K_n . Fix $\epsilon > 0$. We will show that for all sufficiently large n , $K_n < \epsilon$. Let $\eta_\epsilon(x)$ be a bounded continuous function with bounded support such that $\mathbf{E}\{|\eta_\epsilon(X) - \eta(X)|^2\} \leq \frac{\epsilon}{12}$. Since $\mathbf{E}\{Y^2\} < \infty$ implies that $\eta(x) \in L^2(\mathbf{P}_X)$, such a function is assured to exist; the set of bounded

continuous functions with bounded support is dense in $L^2(\mu)$ for all probability measures μ . By (17),

$$\begin{aligned}
K_n &\leq 4\mathbf{E}\left\{\left|2c_n\left(\frac{\sum_{i \in I_V} \bar{\eta}_n(X_i)}{|I_V|} - \frac{1}{2}\right) - \frac{\sum_{i \in I_V} \eta_\epsilon(X_i)}{|I_V|}\right|^2\right\} \\
&\quad + 4\mathbf{E}\left\{\left|\frac{\sum_{i \in I_V} \eta_\epsilon(X_i)}{|I_V|} - \frac{\sum_{i \in I_V} \eta_\epsilon(X)}{|I_V|}\right|^2\right\} \\
&\quad + 4\mathbf{E}\left\{\left|\frac{\sum_{i \in I_V} \eta_\epsilon(X)}{|I_V|} - \eta_\epsilon(X)\right|^2\right\} \\
&\quad + 4\mathbf{E}\{|\eta_\epsilon(X) - \eta(X)|^2\} \\
&\triangleq 4(K_{n1} + K_{n2} + K_{n3} + K_{n4}).
\end{aligned}$$

First, consider K_{n1} .

$$\begin{aligned}
K_{n1} &= \mathbf{E}\left\{\left|\frac{\sum_{i \in I_V} (2c_n(\bar{\eta}_n(X_i) - \frac{1}{2}) - \eta_\epsilon(X_i))}{|I_V|} 1_{\{|I_V|>0\}} - c_n 1_{\{|I_V|=0\}}\right|^2\right\} \\
&\leq 2\mathbf{E}\left\{\left|\frac{\sum_{i \in I_V} (2c_n(\bar{\eta}_n(X_i) - \frac{1}{2}) - \eta_\epsilon(X_i))}{|I_V|} 1_{\{|I_V|>0\}}\right|^2\right\} \\
&\quad + 2\mathbf{E}\{c_n^2 1_{\{|I_V|=0\}}\},
\end{aligned}$$

with the equality following from algebra and the inequality from (17). Then, noting that $|I_V| = \sum_{i=1}^n 1_{\{X_i \in B_{r_n}(X)\}}$ is binomial with parameter $P_{X_1}\{X_1 \in B_{r_n}(X)\}$ when conditioned on X , we continue,

$$\begin{aligned}
K_{n1} &\leq 2\mathbf{E}\left\{\left|\frac{\sum_{i \in I_V} (2c_n(\bar{\eta}_n(X_i) - \frac{1}{2}) - \eta_\epsilon(X_i))}{|I_V|}\right|^2\right\} + 2\mathbf{E}\left\{c_n^2 \left(1 - \mathbf{P}_{X_1}\{X_1 \in B_{r_n}(X)\}\right)^n\right\} \\
&\leq 2c\mathbf{E}\left\{\left|2c_n(\bar{\eta}_n(X) - \frac{1}{2}) - \eta_\epsilon(X)\right|^2\right\} + 2\mathbf{E}\left\{\frac{2c_n^2}{n\mathbf{P}_{X_1}\{X_1 \in B_{r_n}(X)\}}\right\}.
\end{aligned}$$

Here, the second inequality follows for some constant c , in part by applying technical Lemma 5 and in part by noting $(1-x)^n \leq \exp(-nx) \leq \frac{1}{nx}$ for $0 \leq x \leq 1$ and $n = 1, 2, \dots$. Continuing by applying (17), we have

$$K_{n1} \leq 2c\mathbf{E}\left\{\left|2c_n(\bar{\eta}_n(X) - \frac{1}{2}) - \eta(X)\right|^2\right\} + \mathbf{E}\{|\eta_\epsilon(X) - \eta(X)|^2\} + \mathbf{E}\left\{\frac{4c_n^2}{n\mathbf{P}_{X_1}\{X_1 \in B_{r_n}(X)\}}\right\}.$$

For our specific choice of agent decision rules, note that $\bar{\eta}_n(x) = \mathbf{E}\{T_n(Y) | X = x\} = \mathbf{E}\left\{\left(\frac{1}{2c_n}Y + \frac{1}{2}\right)1_{\{|Y| \leq c_n\}} + \frac{1}{2}1_{\{|Y| > c_n\}} \mid X = x\right\}$. Substituting this above and applying Jensen's inequality, we have

$$\begin{aligned}
K_{n1} &\leq 2c\mathbf{E}\left\{\left|\mathbf{E}\{Y 1_{\{|Y| > c_n\}} | X\}\right|^2\right\} + \frac{\epsilon}{12} + \mathbf{E}\left\{\frac{4c_n^2}{n\mathbf{P}_{X_1}\{X_1 \in B_{r_n}(X)\}}\right\} \\
&\leq 2c\mathbf{E}\left\{\mathbf{E}\{Y^2 1_{\{|Y| > c_n\}} | X\}\right\} + \frac{\epsilon}{12} + \mathbf{E}\left\{\frac{4c_n^2}{n\mathbf{P}_{X_1}\{X_1 \in B_{r_n}(X)\}}\right\} \\
&= 2c\mathbf{E}\{Y^2 1_{\{|Y| > c_n\}}\} + \frac{\epsilon}{12} + \mathbf{E}\left\{\frac{4c_n^2}{n\mathbf{P}_{X_1}\{X_1 \in B_{r_n}(X)\}}\right\}. \tag{19}
\end{aligned}$$

Since $f_n(y) = y^2 1_{\{|y| > c_n\}}$ is a monotonically decreasing sequence of functions and $f_n(y) \rightarrow 0$ everywhere, then by the Monotone Convergence Theorem, the first term in (19) converges to zero. The third term in (19) converges to zero by the same argument that was applied for J_n . Thus, $\limsup_{n \rightarrow \infty} K_{n1} \leq \frac{\epsilon}{12}$.

Observe that η_ϵ is uniformly continuous, since by construction, it is a bounded continuous function with bounded support. Let $\delta > 0$ be such that if $\|x - x'\| < \delta$, then $|\eta_\epsilon(x) - \eta_\epsilon(x')| \leq \sqrt{\frac{\epsilon}{12}}$. Since $r_n \rightarrow 0$, for all sufficiently large n , $r_n < \delta$. Thus, for all sufficiently large n ,

$$\begin{aligned} K_{n2} &= \mathbf{E}\left\{\left|\frac{\sum_{i \in I_V} (\eta_\epsilon(X_i) - \eta_\epsilon(X))}{|I_V|}\right|^2\right\} \\ &\leq \frac{\epsilon}{12}, \end{aligned}$$

since for all $i \in I_V$, $\|X_i - X\| \leq r_n$. Next, consider K_{n3} . We have

$$\begin{aligned} K_{n3} &= \mathbf{E}\{\eta_\epsilon(X)^2 1_{\{|I_V|=0\}}\} \\ &\leq \sup_x (\eta_\epsilon(x)^2) \mathbf{E}\{1_{\{|I_V|=0\}}\} \\ &\leq \sup_x (\eta_\epsilon(x)^2) \mathbf{E}\left\{\frac{2c_n^2}{n \mathbf{P}_{X_1}\{X_1 \in B_{r_n}(X)\}}\right\}, \end{aligned}$$

in the usual way, as we see that $K_{n3} \rightarrow 0$. Finally, $K_{n4} \leq \frac{\epsilon}{12}$ by our choice of $\eta_\epsilon(x)$. Thus,

$$\begin{aligned} \limsup_{n \rightarrow \infty} K_n &\leq 4\left(\frac{\epsilon}{12} + \frac{\epsilon}{12} + 0 + \frac{\epsilon}{12}\right) \\ &= \epsilon. \end{aligned}$$

Since ϵ was arbitrary, it is clear that K_n converges to zero. This completes the proof. \blacksquare

VI. DISTRIBUTED REGRESSION WITHOUT ABSTENTION

Finally, let us consider the model for distributed regression without abstention. Now, $\mathcal{Y} = \mathbb{R}$; agents will receive real-valued training data labels Y_i . However, when asked to respond with information, they will reply with either 0 or 1, as abstention is not an option.

In this section, we first establish natural regularity conditions for candidate fusion rules and specify a reasonable class of agent decision rules. As an important negative result, we then demonstrate that for any agent decision rule within this class, there does not exist a regular fusion rule that is L_2 consistent for every distribution \mathbf{P}_{XY} . This result establishes the impossibility of universal consistency in this model for distributed regression without abstention for a restricted, but reasonable class of decision rules.

To begin, consider the set of agent decision rules specified according to (1) for some $\bar{\delta}_n(\cdot)$. In this model without abstention, we require that the implicit responses satisfy $\mathbf{P}\{\delta_{ni} = \text{abstain}\} = 0$, but we impose no additional constraints on the agent decision rules. With the formalism introduced in Section II, this assumption is equivalent to assuming $\{\bar{\delta}_n(\cdot)\}_{n=1}^\infty \subset \mathcal{A} = \{\delta : \mathcal{X} \times \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]\}$.

A fusion rule consists of a sequence of functions $\{\hat{\eta}_n\}_{n=1}^\infty$ mapping $\mathcal{X} \times \mathcal{S}^n$ to $\mathcal{Y} = \mathbb{R}$. Recall from Section II, we can regard $\mathcal{S} = \{1, 0\}$ in this model without abstention. To proceed, we require some regularity on $\{\hat{\eta}_n(\cdot)\}_{n=1}^\infty$. Namely, let us consider only fusion rules that satisfy the following assumptions:

- (A1) $\hat{\eta}_n(x, \cdot)$ is permutation invariant for all $x \in \mathcal{X}$. That is, for all $x \in \mathcal{X}$, any $b \in \{0, 1\}^n$, and any permutation of b , $b' \in \{0, 1\}^n$, $\hat{\eta}_n(x, b) = \hat{\eta}_n(x, b')$.
- (A2) For every $x \in \mathcal{X}$, $\hat{\eta}_n(x, \cdot)$ is Lipschitz in the average Hamming distance. That is, there exists a constant C such that

$$|\hat{\eta}_n(x, b_1) - \hat{\eta}_n(x, b_2)| \leq C \frac{1}{n} \sum_{i=1}^n |b_{1i} - b_{2i}| \quad (20)$$

for every $b_1, b_2 \in \{0, 1\}^n$.

Once again, we will consider $L_n = \mathbf{E}\{|\hat{\eta}_n(X) - Y|^2\}$ with the expectation being taken over X , $D_n = \{(X_i, Y_i)\}_{i=1}^n$, and the randomness introduced in the agent decision rules.

A. Main Result and Comments

The following provides a negative result.

Theorem 4: For every sequence of agent decision rules specified according to (1) with a pointwise convergent sequence of functions $\{\bar{\delta}_n(\cdot)\}_{n=1}^\infty \subset \mathcal{A}$, there is no fusion rule $\{\hat{\eta}_n(\cdot)\}_{n=1}^\infty$ satisfying assumptions (A1) and (A2) such that

$$\lim_{n \rightarrow \infty} \mathbf{E}\{L_n\} = L^* \quad (21)$$

for every distribution \mathbf{P}_{XY} satisfying $\mathbf{E}\{Y^2\} < \infty$.

Note that there is nothing particularly special about the one bit regime and regression. In fact, under the conditions of the theorem, universal consistency cannot be achieved in a multi-class classification problem with even three possible labels. However, we consider regression as it illustrates the point nicely.

The restriction to distributions satisfying $\mathbf{E}\{Y^2\} < \infty$ actually strengthens this negative result, for without such a condition, Theorem 4 is trivial. In the proof, a counter-example is derived where Y is binary-valued, a much stronger case that also satisfies this condition.

Further, the requirement that $\{\bar{\delta}_n(\cdot)\}_{n=1}^\infty$ be pointwise convergent is mild and is only a technical point in the proof. Indeed, the result can be trivially extended to allow for weaker notions of convergence.

B. Proof of Theorem 4

The proof will proceed by specifying two random variables (X, Y) and (X', Y') with $\eta(x) = \mathbf{E}\{Y | X = x\} \neq \mathbf{E}\{Y' | X' = x\} = \eta'(x)$. Asymptotically, however, the fusion center's estimate will be indifferent to whether the agents are trained with random data distributed according to \mathbf{P}_{XY} or $\mathbf{P}_{X'Y'}$. This observation will contradict universal consistency and complete the proof.

Proof: To start, fix a pointwise convergent sequence of functions $\{\bar{\delta}_n(\cdot)\}_{n=1}^\infty \subseteq \mathcal{A}$, arbitrary $x_0, x_1 \in \mathcal{X}$, and distinct $y_0, y_1 \in \mathbb{R}$. Let us specify a distribution \mathbf{P}_{XY} . Let $\mathbf{P}_X\{x_0\} = q$, $\mathbf{P}_X\{x_1\} = 1 - q$, and $\mathbf{P}_{Y|X}\{Y = y_i | X = x_i\} = 1$ for $i = 0, 1$. Clearly, for this distribution $\eta(x_i) = y_i$ for $i = 0, 1$.

Suppose that the ensemble is trained with random data distributed according to \mathbf{P}_{XY} and that the fusion center wishes to classify $X = x_0$. According to the model, after broadcasting X to the agents, the fusion center will observe a random sequence of n bits $\{\delta_{ni}\}_{i=1}^n$. For all $i \in \{1, \dots, n\}$ and all n ,

$$\mathbf{P}\{\delta_{ni} = 1 | X = x_0\} = \bar{\delta}_n(x_0, x_0, y_0)q + \bar{\delta}_n(x_0, x_1, y_1)(1 - q). \quad (22)$$

Now, let us define a sequence of auxiliary random variables, $\{(X'_n, Y')\}_{n=1}^\infty$, with distributions satisfying

$$\begin{aligned} \mathbf{P}_{X'_n}\{x_1\} &= \frac{\bar{\delta}_n(x_0, x_0, y_0)q + \bar{\delta}_n(x_0, x_1, y_1)(1 - q) - \bar{\delta}_n(x_0, x_1, y_1)}{\bar{\delta}_n(x_0, x_0, y_1) - \bar{\delta}_n(x_0, x_1, y_0)} \\ \mathbf{P}_{X'_n}\{x_0\} &= 1 - \mathbf{P}_{X'_n}\{x_1\} \end{aligned}$$

$$\mathbf{P}_{Y'|X'_n}\{Y' = y_{1-i} | X'_n = x_i\} = 1, \quad i = 0, 1.$$

Here, $\eta'(x_i) = \mathbf{E}\{Y' | X'_n = x_i\} = y_{1-i}$. Suppose that the ensemble were trained with random data distributed according to $\mathbf{P}_{X'_n Y'}$ and let $\{\delta_{ni}^{(n)}\}_{i=1}^n$ denote the random response variables of the agents.

Then, we have

$$\begin{aligned} &\mathbf{P}\{\delta_{ni}^{(n)} = 1 | X'_n = x_0\} \\ &= \bar{\delta}_n(x_0, x_0, y_1) \frac{\bar{\delta}_n(x_0, x_0, y_0)q + \bar{\delta}_n(x_0, x_1, y_1)(1 - q) - \bar{\delta}_n(x_0, x_1, y_1)}{\bar{\delta}_n(x_0, x_0, y_1) - \bar{\delta}_n(x_0, x_1, y_0)} \\ &\quad + \bar{\delta}_n(x_0, x_1, y_0) \left(1 - \frac{\bar{\delta}_n(x_0, x_0, y_0)q + \bar{\delta}_n(x_0, x_1, y_1)(1 - q) - \bar{\delta}_n(x_0, x_1, y_1)}{\bar{\delta}_n(x_0, x_0, y_1) - \bar{\delta}_n(x_0, x_1, y_0)}\right) \\ &= \mathbf{P}\{\delta_{ni} = 1 | X = x_0\}, \end{aligned} \quad (23)$$

for all n . Thus, conditioned on the observation to be labeled by the ensemble X (or X'_n), the fusion center will observe an identical stochastic process regardless of whether the ensemble was trained with data distributed according to \mathbf{P}_{XY} or $\mathbf{P}_{X'_n Y'}$ for any fixed n . Note, this observation is true despite the fact that $\eta(x) \neq \eta'(x)$.

Finally, let (X', Y') be such that

$$\mathbf{P}_{X'}\{x_1\} = \lim_{n \rightarrow \infty} \mathbf{P}_{X'_n}\{x_1\} \quad (24)$$

$$\mathbf{P}_{X'}\{x_0\} = 1 - \mathbf{P}_{X'}\{x_1\}$$

$$\mathbf{P}_{Y'|X'}\{Y' = y_{1-i} | X' = x_i\} = 1, \quad i = 0, 1.$$

Again, $\eta'(x_i) = \mathbf{E}\{Y' | X' = x_i\} = y_{1-i}$. These limits are assured to exist by the assumption that $\{\bar{\delta}_n(\cdot)\}_{n=1}^\infty$ is a pointwise converging sequence of functions. Finally, let $\{\delta'_{ni}\}_{i=1}^n$ denote the random response random variables for the ensemble agents trained with data distributed according to $\mathbf{P}_{X'Y'}$.

By standard orthogonality arguments [12], for the ensemble to be universally consistent, we must have both

$$\mathbf{E}\{|\hat{\eta}_n(X, \{\delta_{ni}\}_{i=1}^n) - \eta(X)|^2\} \rightarrow 0 \quad (25)$$

and

$$\mathbf{E}\{|\hat{\eta}_n(X', \{\delta'_{ni}\}_{i=1}^n) - \eta'(X')|^2\} \rightarrow 0. \quad (26)$$

Let us assume that (25) holds; we now demonstrate that necessarily,

$$\mathbf{E}\{|\hat{\eta}_n(X', \{\delta'_{ni}\}_{i=1}^n) - \eta(X')|^2\} \rightarrow 0. \quad (27)$$

Since $\eta(x) \neq \eta'(x)$, (27) contradicts (26) and the proposition of universal consistency. To show (27), it suffices to focus on the L^2 risk conditioned on X' , due to the convenient point-mass structure of $\mathbf{P}_{X'}$. To proceed, note that by (17), for any $b \in \{0, 1\}^n$,

$$\begin{aligned} \mathbf{E}\{|\hat{\eta}_n(X', \{\delta'_{ni}\}_{i=1}^n) - \eta(X')|^2 | X' = x_0\} &\leq 2\mathbf{E}\{|\hat{\eta}_n(X', b) - \eta(X')|^2 | X' = x_0\} \\ &\quad + 2\mathbf{E}\{|\hat{\eta}_n(X', \{\delta'_{ni}\}_{i=1}^n) - \hat{\eta}_n(X', b)|^2 | X' = x_0\} \\ &\triangleq 2T_1(b) + 2T_2(b). \end{aligned}$$

In particular, let us select $b \in \{0, 1\}^n$ randomly such that the components are i.i.d. with $b_i \sim \mathbf{P}\{\delta_{ni} | X = x_0\}$ for all $i = 1, \dots, n$. Note that if we can show that $\mathbf{E}_b\{T_1(b) + T_2(b)\} \rightarrow 0$, then the result holds by the probabilistic method. First consider $T_1(b)$. Note that we have

$$\begin{aligned} \mathbf{E}_b\{T_1(b)\} &= \mathbf{E}\{|\hat{\eta}_n(X', b) - \eta(X')|^2 | X' = x_0\} \\ &= \mathbf{E}\{|\hat{\eta}_n(X, \{\delta_{ni}\}_{i=1}^n) - \eta(X)|^2 | X = x_0\}, \end{aligned}$$

by our selection of b . Thus, $\mathbf{E}_b\{T_1(b)\}$ must converge to zero by the assumption that (25) holds true. Considering $T_2(b)$, note that

$$\begin{aligned} \mathbf{E}_b\{T_2(b)\} &= \mathbf{E}\{|\hat{\eta}_n(X', b) - \hat{\eta}_n(X', \{\delta'_{ni}\}_{i=1}^n)|^2 | X' = x_0\} \\ &\leq C^2 \mathbf{E}\left\{\left|\frac{1}{n} \sum_{i=1}^n b_i - \frac{1}{n} \sum_{i=1}^n \delta'_{ni}\right|^2 \middle| X' = x_0\right\} \\ &\leq 3C^2 \mathbf{E}\left\{\left|\frac{1}{n} \sum_{i=1}^n b_i - \mathbf{P}\{\delta_{ni} = 1 | X = x_0\}\right|^2\right\} \end{aligned} \quad (28)$$

$$+ 3C^2 \mathbf{E}\left\{\left|\frac{1}{n} \sum_{i=1}^n \delta'_{ni} - \mathbf{P}\{\delta'_{ni} = 1 | X' = x_0\}\right|^2 \middle| X' = x_0\right\} \quad (29)$$

$$+ 3C^2 |\mathbf{P}\{\delta_{ni} = 1 | X = x_0\} - \mathbf{P}\{\delta'_{ni} = 1 | X' = x_0\}|^2. \quad (30)$$

Here, the first inequality follows from assumptions (A1) and (A2) and the second inequality follows by (17). Note that since $\{b_i\}_{i=1}^n$ is i.i.d. with $b_i \sim \mathbf{P}\{\delta_{ni} = 1 | X = x_0\}$,

$$3C^2 \mathbf{E}\left\{\left|\frac{1}{n} \sum_{i=1}^n b_i - \mathbf{P}\{\delta_{ni} = 1 | X = x_0\}\right|^2\right\} \leq \frac{3C^2}{4n},$$

after bounding the variance of a binomial random variable; therefore, (28) must converge to zero. A similar argument can be applied to (29). Next, from (23),

$$|\mathbf{P}\{\delta_{ni} = 1 | X = x_0\} - \mathbf{P}\{\delta'_{ni} = 1 | X' = x_0\}|^2 = |\mathbf{P}\{\delta_{ni}^{(n)} = 1 | X'_n = x_0\} - \mathbf{P}\{\delta'_{ni} = 1 | X' = x_0\}|^2.$$

Thus, (30) must converge to zero by our design of (X', Y') in (24). Finally, we have demonstrated that (27) holds true; by the discussion above, this completes the proof. \blacksquare

VII. CONCLUSIONS AND FUTURE WORK

Motivated by sensor networks and other distributed settings, this paper has presented several models for distributed learning. The models differ from classical works in statistical pattern recognition by allocating observations of an i.i.d. sampling process to individual learning agents. By limiting the ability of the agents to communicate, we constrain the amount of information available to the ensemble and to the fusion center for use in classification or regression. This setting models a distributed environment and presents new questions to consider with regard to universal consistency.

Insofar as these models present a useful picture of distributed scenarios, this paper has answered several questions about whether or not the guarantees provided by Stone's Theorem in centralized environments hold in distributed settings. The models have demonstrated that when agents are allowed to communicate $\log_2(3)$ bits per decision, the ensemble can achieve universal consistency in both binary classification

and regression frameworks in the limit as the number of agents increases without bound. In the binary classification case, we have demonstrated this property as a special case of naive kernel classifiers. In the regression case, we have shown this to hold true with randomized agent decision rules. When investigating the necessity of these $\log_2(3)$ bits, we have found that in the binary classification framework only one bit per agent per classification was necessary for universal consistency, and the analysis provided an interesting comparison for naive kernel methods in the traditional framework. For regression, we have established the impossibility of universal consistency in the one bit regime for a natural, but restricted class of candidate rules.

With regard to future research in distributed learning, there are numerous directions of interest. As these results are useful only if they accurately depict some aspect of distributed environments, other perhaps more reflective models are important to consider. In particular, the current models assume that a reliable physical layer exists where bits transmitted from the agents are guaranteed to arrive unperturbed at the fusion center. Future research may consider richer model for this communication, perhaps within an information-theoretic (i.e., Shannon-theoretic) formalism. Further, the current models consider simplified network models where the fusion center communicates with agents via a broadcast medium and each agent has a direct, albeit limited, channel to the fusion center. Future research may focus on network models that allow for inter-agent communication. Consistent with the spirit of sensor networks, we might allow agents to communicate locally amongst themselves (or perhaps, hierarchically) before coordinating a response to the fusion center. In general, models of this form would weaken (A) in the discussion in Section II by allowing for correlated agent responses. A related assumption in this work is that the underlying data is i.i.d. Extending the results to other sampling process is important since in many distributed applications, the data observed by the agents may be correlated. In this vein, connections to results in statistical pattern recognition results under non-i.i.d. sampling processes would be interesting and important to resolve.

Finally, from a learning perspective, the questions we have considered in this paper have been focused on the statistical issue of universal consistency. Though such a consideration seems to be one natural first step, other comparisons between centralized and distributed learning are essential, perhaps with respect to convergence rate and the finite data reality that exists in any practical system. Such questions open the door for agents to receive multiple training examples and may demand more complicated local decision algorithms; in particular, it may be interesting to study local regularization strategies for agents in an ensemble. Future work may explore these and other questions frequently explored in traditional, centralized learning systems, with the hope of further understanding the nature of distributed learning

under communication constraints.

APPENDIX

This appendix includes important facts that are commonly used in the study of nonparametric statistics and are similarly applied in the proofs above. Lemma 1 is a basic result from probability theory and is included for clarity. Lemma 2 follows from Theorem 23.2 and Lemma 23.6 in [12] applied to the naive kernel. The proof of Theorem 6.2 in [7] contains the fundamental steps needed to prove Lemma 3. Lemma 4 can be found as Lemma 4.1 in [12]. Lemma 5 follows from arguments used in proving Theorem 5.1 in [12] applied to the naive kernel.

Lemma 1: Suppose $\{X_n\}_{n=1}^\infty$ is a sequence of random variables such that $X_n \rightarrow X$ in probability. Then, for any sequence of events $\{A_n\}_{n=1}^\infty$ with $\liminf \mathbf{P}\{A_n\} > 0$,

$$\mathbf{P}\{|X_n - X| > \epsilon | A_n\} \rightarrow 0.$$

for all $\epsilon > 0$.

Proof: After noting that,

$$\begin{aligned} \mathbf{P}\{|X_n - X| > \epsilon\} &= \mathbf{P}\{|X_n - X| > \epsilon | A_n\} \mathbf{P}\{A_n\} + \mathbf{P}\{|X_n - X| > \epsilon | \bar{A}_n\} \mathbf{P}\{\bar{A}_n\} \\ &\geq \mathbf{P}\{|X_n - X| > \epsilon | A_n\} \mathbf{P}\{A_n\}, \end{aligned}$$

the Lemma follows trivially from the fact that $\liminf \mathbf{P}\{A_n\} > 0$ and $X_n \rightarrow X$ in probability. The proof follows similarly if $X_n \rightarrow \infty$ in probability. ■

Lemma 2: Let $X \sim \mathbf{P}_X$ be an \mathbb{R}^d -valued random variable and fix any function $f \in L(\mathbf{P}_X)$. For an arbitrary sequence of real numbers $\{r_n\}_{n=1}^\infty$, define a sequence of functions $f_n(x) = \mathbf{E}\{f(X) | X \in B_{r_n}(x)\}$. If $r_n \rightarrow 0$, then $f_n(X) \rightarrow f(X)$ in probability.

Lemma 3: Let $X \sim \mathbf{P}_X$ be an \mathbb{R}^d -valued random variable and define $\{r_n\}_{n=1}^\infty$ and $\{a_n\}_{n=1}^\infty$ as arbitrary sequences of real numbers such that $r_n \rightarrow 0$ and $a_n \rightarrow \infty$. If $(r_n)^d a_n \rightarrow \infty$, then

$$a_n \int 1_{B_{r_n}(X)}(y) P_X(dy) \rightarrow \infty \text{ i.p.}$$

Lemma 4: Suppose $B(n, p)$ is a binomially distributed random variable with parameters n and p . Then,

$$\mathbf{E}\left\{\frac{1}{B(n, p)} 1_{\{B(n, p) > 0\}}\right\} \leq \frac{2}{(n+1)p}.$$

Lemma 5: There is a constant c such that for any measurable function f , any \mathbb{R}^d -valued random variable X , and any sequence $\{r_n\}_{n=1}^\infty$,

$$\mathbf{E}\left\{\frac{\sum_{i=1}^n 1_{\{X_i \in B_{r_n}(X)\}} f(X_i)}{\sum_{i=1}^n 1_{\{X_i \in B_{r_n}(X)\}}}\right\} \leq c \mathbf{E}\{f(X)\}$$

for all n .

REFERENCES

- [1] I. F. Akyildiz, W. Su, Y. Sankarasubramaniam, and E. Cayirci, "A survey on sensor networks," *IEEE Communications Magazine*, vol. 40, no. 8, pp. 102–114, 2002.
- [2] M. Barkat and P. K. Varshney, "Decentralized CFAR signal detection," *IEEE Trans. Aerospace and Electronic Systems*, vol. 25, Mar. 1989.
- [3] R. Blum, S. Kassam, and H. V. Poor, "Distributed detection with multiple sensors II: Advanced topics," *Proceedings of the IEEE*, vol. 85, pp. 64–79, Jan. 1997.
- [4] L. Breiman, "Bagging predictors," *Machine Learning*, vol. 26, no. 2, pp. 123–140, 1996.
- [5] T. M. Cover, "Rates of convergence for nearest neighbor procedures," *Proc. 1st Annu. Hawaii Conf. Systems Theory*, pp. 413–415, Jan. 1968.
- [6] A. D'Costa and A. M. Sayeed, "Collaborative signal processing for distributed classification in sensor networks," in *Lecture Notes in Computer Science (Proceedings of IPSN'03)*, F. Zhao and L. Guibas, Eds. Berlin: Springer, 2003.
- [7] L. Devroye, L. Györfi, and G. Lugosi, *A Probabilistic Theory of Pattern Recognition*. New York: Springer, 1996.
- [8] M. Dong, L. Tong, and B. M. Sadler, "Information retrieval and processing in sensor networks: Deterministic scheduling vs. random access," submitted to *IEEE Trans. on Inform. Theory*, 2004.
- [9] Y. Freund and R. E. Schapire, "A decision-theoretic generalization of on-line learning and an application to boosting," *Computer and System Sciences*, vol. 55, pp. 119–139, 1997.
- [10] Y. Freund, R. E. Schapire, Y. Singer, and M. K. Warmuth, "Using and combining predictors that specialize," in *Proceedings of the Twenty-Ninth Annual ACM Symposium on the Theory of Computing*, El Paso, Texas, 1997, pp. 334–343.
- [11] W. Greblicki and M. Pawlak, "Necessary and sufficient conditions for Bayes risk consistency of recursive kernel classification rule," *IEEE Trans. Inform. Theory*, vol. IT-33, pp. 408–412, 1987.
- [12] L. Györfi, M. Kohler, A. Krzyżak, and H. Walk, *A Distribution-Free Theory of Nonparametric Regression*. New York: Springer, 2002.
- [13] R. Jacobs, M. I. Jordan, S. Nowlan, and G. E. Hinton, "Adaptive mixtures of local experts," *Neural Computation*, vol. 3, no. 1, pp. 125–130, 1991.
- [14] M. Kearns and H. S. Seung, "Learning from a population of hypotheses," *Machine Learning*, vol. 18, pp. 255–276, 1995.
- [15] J. Kittler, M. Hatef, P. W. Duin, and J. Matas, "On combining classifiers," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 20, no. 3, pp. 226–239, 1998.
- [16] A. N. Kolmogorov and S. V. Fomin, *Introductory Real Analysis*. New York: Dover, 1975.
- [17] J. H. Kotecha, V. Ramachandran, and A. Sayeed, "Distributed multi-target classification in wireless sensor networks," to appear in *IEEE Journal of Selected Areas in Communications (Special Issue on Self-Organizing Distributed Collaborative Sensor Networks)*, July 2003.
- [18] A. Krzyżak, "The rates of convergence of kernel regression estimates and classification rules," *IEEE Trans. Inform. Theory*, vol. IT-32, pp. 668–679, 1986.
- [19] S. R. Kulkarni and S. E. Posner, "Rates of convergence of nearest neighbor estimation under arbitrary sampling," *IEEE Trans. Inform. Theory*, vol. 41, no. 4, pp. 1028–1039, July 1995.
- [20] S. R. Kulkarni, S. E. Posner, and S. Sandilya, "Data-dependent k_n -nn and kernel estimators consistent for arbitrary processes," *IEEE Trans. Inform. Theory*, vol. 48, no. 10, pp. 2785–2788, 2002.
- [21] S. Kumar, F. Zhao, and D. Shephard, "Collaborative signal and information processing in microsensors networks," *IEEE Signal Processing Magazine*, vol. 19, no. 2, pp. 13–14, 2002.

- [22] A. Lazarevic and Z. Obradovic, "The distributed boosting algorithm," in *KDD '01: Proceedings of the Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA*. ACM Press, 2001, pp. 311–316.
- [23] D. Li, K. Wong, Y. H. Hu, and A. Sayeed, "Detection, classification, and tracking of targets," *IEEE Signal Processing Magazine*, vol. 19, no. 2, pp. 17–29, 2002.
- [24] G. Lugosi, "Learning with an unreliable teacher," *Pattern Recognition*, vol. 25, pp. 79–87, 1992.
- [25] G. Morvai, S. R. Kulkarni, and A. B. Nobel, "Regression estimation from an individual stable sequence," *Statistics*, vol. 33, pp. 99–119, 1999.
- [26] X. Nguyen, M. J. Wainwright, and M. I. Jordan, "Decentralized detection and classification using kernel methods," in *Proceedings of the Twenty-first International Conference on Machine Learning, Banff, Canada, 2004*.
- [27] A. B. Nobel, "Limits to classification and regression estimation from ergodic processes," *Ann. Statist.*, vol. 27, pp. 262–273, 1999.
- [28] A. B. Nobel and T. M. Adams, "On regression estimation from ergodic samples with additive noise," *IEEE Trans. Inform. Theory*, vol. 47, pp. 2895–2902, 2001.
- [29] A. B. Nobel, G. Morvai, and S. Kulkarni, "Density estimation from an individual sequence," *IEEE Trans. Inform. Theory*, vol. 44, pp. 537–541, Mar. 1998.
- [30] J. B. Predd, S. R. Kulkarni, and H. V. Poor, "Consistency in a model for distributed learning with specialists," in *Proceedings of the 2004 IEEE International Symposium on Information Theory, Chicago, IL, June 2004*.
- [31] —, "Consistency in models for communication constrained distributed learning," in *Learning Theory, Proceedings of the 17th Annual Conference on Learning Theory, COLT 2004, Banff, Canada*, J. Shawe-Taylor and Y. Singer, Eds. Springer, July 2004.
- [32] —, "Distributed learning in wireless sensor networks," in *Proceedings of the 42nd Annual Allerton Conference on Communication, Control, and Computing, Monticello, IL, Sept. 2004*.
- [33] G. Roussas, "Nonparametric estimation in markov processes," *Ann. Inst. Statist. Math.*, vol. 21, pp. 73–87, 1967.
- [34] S. N. Simic, "A learning theory approach to sensor networks," *IEEE Pervasive Computing*, vol. 2, no. 4, pp. 44–49, 2003.
- [35] C. J. Stone, "Consistent nonparametric regression," *Annals of Statistics*, vol. 5, pp. 595–645, 1977.
- [36] J. N. Tsitsiklis, "Decentralized detection," in *Advances in Statistical Signal Processing*. JAI Press, 1993, pp. 297–344.
- [37] P. K. Varshney, *Distributed Detection and Data Fusion*. New York: Springer, 1996.
- [38] V. V. Veeravalli, "Decentralized quickest change detection," *IEEE Trans. on Inform. Theory*, vol. 47, no. 4, pp. 1657–1656, 2001.
- [39] S. Yakowitz, "Nonparametric density and regression estimation from Markov sequences without mixing assumptions," *J. Multivar. Anal.*, vol. 30, pp. 124–136, 1989.
- [40] —, "Nearest neighbor regression estimation for null-recurrent Markov time series," *Stoch. Processes Appl.*, vol. 48, pp. 311–318, 1993.