

Distributed Kernel Regression: An Algorithm for Training Collaboratively

J. B. Predd, *Member, IEEE*, S. R. Kulkarni, *Fellow, IEEE*, and H. V. Poor, *Fellow, IEEE*

Abstract—This paper addresses the problem of distributed learning under communication constraints, motivated by distributed signal processing in wireless sensor networks and data mining with distributed databases. After formalizing a general model for distributed learning, an algorithm for collaboratively training regularized kernel least-squares regression estimators is derived. Noting that the algorithm can be viewed as an application of successive orthogonal projection algorithms, its convergence properties are investigated and the statistical behavior of the estimator is discussed in a simplified theoretical setting.

I. INTRODUCTION

In this paper, we address the problem of *distributed learning under communication constraints*, motivated primarily by distributed signal processing in wireless sensor networks (WSNs) and data mining with distributed databases. WSNs are *a fortiori* designed to make inferences from the environments they are sensing; however they are typically characterized by constraints on energy and bandwidth, which limit the sensors' ability to share data with each other or with a centralized fusion center. In data mining with distributed databases, multiple agents (e.g., corporations) have access to possibly overlapping databases, and wish to collaborate to make optimal inferences; privacy or security concerns, however, may preclude them from fully sharing information. Nonparametric methods studied within machine learning have demonstrated widespread empirical success in many centralized (i.e., communication *unconstrained*) signal processing applications. Thus, in both the aforementioned applications, a natural question arises: can the power of machine learning methods be tapped for nonparametric inference in distributed learning under communication constraints?

In this paper, we address this question by formalizing a general model for distributed learning, and then deriving a distributed algorithm for collaborative training in regularized kernel least-squares regression. The algorithm can be viewed as an instantiation of successive orthogonal projection algorithms, and thus, insight into the statistical behavior of these algorithms can be gleaned from standard analyses in mathematical programming.

A. Related Work

Distributed learning has been addressed in a variety of other works. Reference [9] considered a PAC-like model for learning with many individually trained hypotheses in a distribution-specific learning framework. Reference [13] considered the classical model for decentralized detection [17] in a nonparametric setting. Reference [15] studied the existence of consistent

estimators in several models for distributed learning. From a data mining perspective, [6] and [12] derived algorithms for distributed boosting. Most similar to the research presented here, [7] presented a general framework for distributed linear regression motivated by WSNs.

Ongoing research in the machine learning community seeks to design statistically sound learning algorithms that scale to large data sets (e.g., [3] and references therein). One approach is to decompose the database into smaller “chunks”, and subsequently parallelize the learning process by assigning distinct processors/agents to each of the chunks. In principle, algorithms for parallelizing learning may be useful for distributed learning, and vice-versa. To our knowledge, there has not been an attempt to parallelize reproducing kernel methods using the approach outlined below.

A related area of research lies in the study of ensemble methods in machine learning; examples of these techniques include bagging, boosting, and mixtures of experts (e.g., [5] and others). Typically, the focus of these works is on the statistical and algorithmic advantages of learning with an ensemble and not on the problem of learning under communication constraints. To our knowledge, the methods derived here have not been derived in this related context, though future work in distributed learning may benefit from the many insights gleaned from this important area.

Those familiar with the online learning framework may find our collaborative training algorithm reminiscent of the equations for additive gradient updates [11]. Though both algorithms may be interpreted in the context of successive orthogonal projection algorithms, it does not appear possible to specialize the current model for distributed learning in a way that recovers the online learning framework (or vice versa).

The research presented here generalizes the model and algorithm discussed in [14], which focused exclusively on the WSN application. Distinctions between the current and former work are discussed in more detail below.

B. Organization

The remainder of this paper is organized as follows. In Section II, we review preliminary background information necessary for the remainder of the work. In Section III, we describe a general model for distributed learning and propose a distributed algorithm for collaboratively training regularized kernel least-squares regression estimators. Subsequently, we analyze the algorithm's convergence properties and use these properties to gain insight into the statistical behavior of the estimator in a simplified setting. We conclude with a discussion of the method in Section IV.

II. PRELIMINARIES

In this section, we briefly review the supervised learning model for nonparametric least-squares regression, reproducing kernel

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 Grants CCR-0020524 and CCR-0312413, and in part by the U. S. Army Pantheon Project.

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

methods, and alternating projection algorithms. Since a thorough introduction to these models and methods is beyond the scope of this paper, we refer the reader to standard references on the topics; see, for example, [4], [8], [16] and references therein.

A. Nonparametric Least-squares Regression

Let X and Y be \mathcal{X} and \mathcal{Y} -valued random variables, respectively. \mathcal{X} is known as the feature, input, or observation space; \mathcal{Y} is known as the label, output, or target space. For now, we allow \mathcal{X} to be arbitrary, but take $\mathcal{Y} = \mathbb{R}$. In the least-squares estimation problem, we seek a decision rule mapping inputs to outputs that minimizes the expected squared error. In particular, we seek a function $g : \mathcal{X} \rightarrow \mathcal{Y}$ that minimizes

$$\mathbf{E}\{|g(X) - Y|^2\}.$$

It is well-known that $\eta(x) = \mathbf{E}\{Y | X = x\}$ is the loss minimizing rule. However, without prior knowledge of the joint distribution of (X, Y) , this regression function cannot be computed. In the supervised learning model, one is instead provided a database $S = \{(x_i, y_i)\}_{i=1}^n$ of training examples with $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} \forall i \in \{1, \dots, n\}$; the learning task is to use S to estimate $\eta(x)$.

B. Regularized Kernel Methods

Regularized kernel methods [16] offer one approach to nonparametric regression. In particular, let \mathcal{H}_K denote the *reproducing kernel Hilbert space* (RKHS) induced by a *positive semi-definite kernel* $K(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$; let $\|\cdot\|_{\mathcal{H}_K}$ denote the norm associated with \mathcal{H}_K . In practice, the kernel K is a design parameter, chosen as a similarity measure between inputs to reflect prior application-specific domain knowledge. The regularized kernel least-squares estimate is defined as the solution $f_\lambda \in \mathcal{H}_K$ of the following optimization problem:

$$\min_{f \in \mathcal{H}_K} \left[\sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}_K}^2 \right]. \quad (1)$$

The statistical behavior of this estimator is well-understood under various assumptions on the stochastic process that generates the examples $\{(x_i, y_i)\}_{i=1}^n$ [16], [19]. In this paper, we focus primarily on algorithmic aspects of computing a solution to (1) (or an approximation thereof) in distributed environments. To this end, consider the following ‘‘Representer Theorem’’ proved originally in [10].

Theorem 1 ([10]): Let $f_\lambda \in \mathcal{H}_K$ be the minimizer of (1). Then, there exists $c_\lambda \in \mathbb{R}^n$ such that

$$f_\lambda(\cdot) = \sum_{i=1}^n c_{\lambda,i} K(\cdot, x_i).$$

From a computational perspective, the result is significant because it states that while the objective function (1) is defined over a potentially infinite dimensional Hilbert space, its minimizer must lie in a finite dimensional subspace.

Finally, note that (1) can be naturally interpreted as an orthogonal projection. In particular, by introducing an auxiliary vector $\mathbf{z} \in \mathbb{R}^n$, (1) can be rewritten as the following optimization program:

$$\min \quad \|\mathbf{z} - \mathbf{y}\|_2^2 + \lambda \|f\|_{\mathcal{H}_K}^2 \quad (2)$$

$$\begin{aligned} \text{s.t.} \quad & z_i = f(x_i) \quad \forall i \in \{1, \dots, n\} \\ & \mathbf{z} \in \mathbb{R}^n \\ & f \in \mathcal{H}_K. \end{aligned} \quad (3)$$

Through the constraints in (3), (1) and (2) are equivalent in the following sense: if f_λ is the minimizer of (1) and $(\mathbf{z}', f'_\lambda)$ is the solution of (2), then $f'_\lambda = f_\lambda$. Therefore, through (2), we can interpret the regularized kernel least-squares estimator as a projection of the vector $(\mathbf{y}, 0) \in \mathbb{R}^n \times \mathcal{H}_K$ onto the set

$$\left\{ (\mathbf{z}, f) \in \mathbb{R}^n \times \mathcal{H}_K : z_i = f(x_i) \forall i \in \{1, \dots, n\} \right\} \subset \mathbb{R}^n \times \mathcal{H}_K.$$

This simple observation will recur in the sequel.

C. Alternating Projections Algorithms

Let \mathcal{X} be a Hilbert space with a norm denoted by $\|\cdot\|$. Let C_1, \dots, C_m be closed convex subsets of \mathcal{X} whose intersection $C = \cap_{i=1}^m C_i$ is nonempty. Let $P_C(\hat{x})$ denote the orthogonal projection of $\hat{x} \in \mathcal{X}$ onto C , i.e.,

$$P_C(\hat{x}) \triangleq \arg \min_{x \in C} \|x - \hat{x}\|.$$

Define $P_{C_i}(\hat{x})$ analogously.

Successive orthogonal projection (SOP) algorithms [4] provide a natural way to compute $P_C(\cdot)$ given $\{P_{C_i}(\cdot)\}_{i=1}^m$. For example, the (unrelaxed) SOP algorithm is defined as follows:

$$x_0 := \hat{x} \quad x_n := P_{C_{(n \bmod m) + 1}}(x_{n-1}). \quad (4)$$

In words, the algorithm successively and iteratively projects onto each of the subsets. In the case where C_i is a linear subspace for all $i \in \{1, \dots, m\}$, this algorithm was first studied by von Neumann [18]. Often examined in the context of the *convex feasibility problem*, SOP has been generalized in various ways [4], to address more general convex sets and non-orthogonal (e.g., Bregman) projections; accordingly, the algorithm often takes on other names (e.g., the von Neumann-Halperin algorithm, Bregman’s algorithm). Much of the behavior of this algorithm can be understood through Theorem 2; the proof of this fundamental result can be found in [2].

Theorem 2: Let $\{C_i\}_{i=1}^m$ be a set of closed, convex subsets of \mathcal{X} whose intersection $C = \cap_{i=1}^m C_i$ is nonempty. Let x_n be defined as in (4). Then, for every $x \in C$ and every $n \geq 1$,

$$\|x_n - x\| \leq \|x_{n-1} - x\|.$$

Moreover, $\lim_{n \rightarrow \infty} x_n \in \cap_{i=1}^m C_i$. If C_i are affine for all $i \in \{1, \dots, m\}$, then $\lim_{n \rightarrow \infty} \|x_n - P_C(\hat{x})\| = 0$.

III. DISTRIBUTED KERNEL REGRESSION

A. The Model

In contrast to the model for supervised learning reviewed in Section II, suppose that each member of a collection of m learning agents has limited access to the training database $S = \{(x_i, y_i)\}_{i=1}^n$. In particular, assume that learning agent i has access only to the training examples in subset $S_i \subseteq S$. For convenience, we shall henceforth refer to $\{S_i\}_{i=1}^m$ as an *ensemble*.

A bipartite graph is a convenient way to represent an ensemble in this model for distributed regression. As depicted in Figure 1, nodes on the top-level of the graph represent learning agents; nodes on the bottom-level represent training examples. An edge between a learning agent i and a training sample j signifies that agent i has access to example j , i.e., $(x_j, y_j) \in S_i$. For now, we make no additional assumptions on the structural relationship between the agents’ locally accessible training sets; for example,

we do not require the ensemble $\{S_i\}_{i=1}^m$ to partition S , nor do we require the corresponding bipartite graph to be connected in any way.

To be concrete, consider a few examples that illustrate special-cases of the general model depicted in Figure 1. The standard centralized model for supervised learning can be represented by the graph in Figure 2, where each of the m learning agents has access to all exemplars in the training database. Figure 3 illustrates an ensemble where a publicly available database is available to all the learning agents, each of which retains a private training set. In some applications, \mathcal{X} may be endowed with a topology. For example, in wireless sensor networks, $\mathcal{X} = \mathbb{R}^2$ may model locations in a city; learning agents (i.e., sensors) may exist as points within \mathcal{X} , and query those examples that are “nearby” with respect to the underlying topology; such an ensemble is depicted in Figure 4.

As mentioned earlier, the current model is a generalization of the the work discussed in [14]. Whereas [14] focuses exclusively on the WSN application by assuming a topology on \mathcal{X} and by modeling one agent per training observation, the present formulation allows a more general structure with multiple agents per training datum and an arbitrary input space.

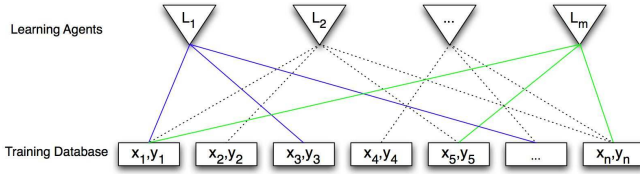


Fig. 1. A Bipartite Graph Representation of an Ensemble in this Model for Distributed Regression

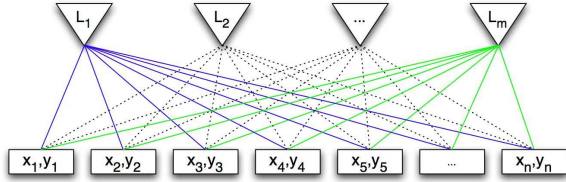


Fig. 2. A “Centralized” Ensemble

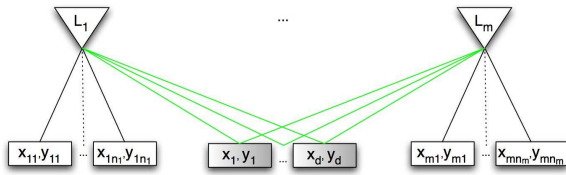


Fig. 3. An Ensemble with a Public Database

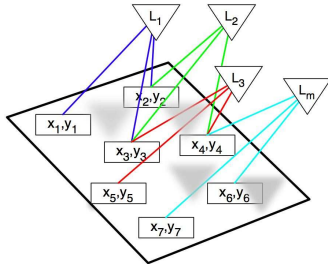


Fig. 4. A Sensor Network: An Ensemble with Topology Dependent Structure

Presumably, each of the m agents wishes to use nonparametric methods to estimate the regression function. One simple approach

is for agent i to compute f_{λ_i} using only the exemplars in its local training database S_i . However, doing so ignores the structure of distributed regression and fails to exploit an opportunity to collaborate using the (partially) shared training database.

We henceforth assume that after locally computing $f_{\lambda_i} \in \mathcal{H}_K$, agent i may share $f_{\lambda_i}(x_j) \in \mathbb{R}$ with any agent k such that $(x_j, y_j) \in S_k$. In other words, neighboring agents (with respect to the bi-partite graph) communicate point estimates for the training data they share. Using such limited communication, can the agents collaborate to jointly improve the accuracy of their estimates?

In the next section, we derive a collaborative training algorithm in this model for distributed nonparametric regression. The algorithm is derived as an application of SOP algorithms applied to a relaxation of the classical regularized kernel least-squares estimator. Subsequently, we analyze its convergence properties and investigate its statistical properties in a simplified theoretical setting.

B. A Collaborative Training Algorithm

For technical convenience, let us introduce sets $\{\bar{S}_i\}_{i=1}^m$, such that $\bar{S}_i \subseteq \{1, \dots, n\}$. Let $j \in \bar{S}_i$ if and only if $(x_j, y_j) \in S_i$. In other words \bar{S}_i contains the indices of the training examples in S_i as enumerated in S . Analogously, let $\bar{S} = \{1, \dots, n\}$.

To begin, let us rewrite (1) in a way that reveals the structure of distributed regression. To do so, first let us introduce a function $f_i \in \mathcal{H}_K$ for each agent $i \in \{1, \dots, m\}$, and consider the following constrained optimization program:

$$\min \quad \|\mathbf{z} - \mathbf{y}\|_2^2 + \sum_{i=1}^m \lambda_i \|f_i\|_{\mathcal{H}_K}^2 \quad (5)$$

$$\text{s.t.} \quad \begin{aligned} z_j &= f_i(x_j) \quad \forall j \in \bar{S}_i, i \in \{1, \dots, m\} \\ f_i &\in \mathcal{H}_K \quad i \in \{1, \dots, m\} \end{aligned} \quad (6)$$

Here, the optimization variables are $\mathbf{z} \in \mathbb{R}^n$ and $\{f_i\}_{i=1}^m \subset \mathcal{H}_K$; $S = \{(x_i, y_i)\}_{i=1}^n$ and $\{\lambda_i\}_{i=1}^m \subset \mathbb{R}$ are the program data. The *coupling constraints* in (6) dictate that for any feasible solution to (5), every agent’s associated function is equivalent when evaluated at $\{x_i\}_{i=1}^n$. As a result, one can think about (5) as an equivalent form of (1) in the following sense.

Lemma 1: Let $(\mathbf{z}, f_{\lambda_1}, \dots, f_{\lambda_m}) \in \mathbb{R}^n \times \mathcal{H}_K^m$ denote the solution of (5) and let $f_{\lambda} \in \mathcal{H}_K$ denote the solution of (1). Assume that $\lambda_i > 0 \forall i \in \{1, \dots, m\}$. Then, $f_{\lambda_1} = \dots = f_{\lambda_m}$. If $\sum_{i=1}^m \lambda_i = \lambda$, then $f_{\lambda} = f_{\lambda_1}$.

This form of the regularized least-squares regression problem suggests a natural relaxation that allows us to incorporate the structure of the distributed regression model into the estimator. In particular, we relax the coupling constraints to require that agents agree only on training examples they share:

$$\min \quad \|\mathbf{z} - \mathbf{y}\|_2^2 + \sum_{i=1}^m \lambda_i \|f_i\|_{\mathcal{H}_K}^2 \quad (7)$$

$$\text{s.t.} \quad \begin{aligned} z_j &= f_i(x_j) \quad \forall j \in \bar{S}_i, i \in \{1, \dots, m\} \\ f_i &\in \mathcal{H}_K \quad i \in \{1, \dots, m\} \end{aligned} \quad (8)$$

Thus, for any feasible solution to (7), $f_i(\mathbf{x}_j) = f_k(\mathbf{x}_j)$ if $(x_j, y_j) \in S_i \cap S_k$. Looked at in this way, (5) models the “centralized ensemble” depicted in Figure 2, while (7) captures the more general structure in Figure 1.

Note that just as (1) can be interpreted as a projection via (2), (7) can be interpreted as a (weighted) projection of the vector

$(\mathbf{y}, 0, \dots, 0) \in \mathbb{R}^n \times \mathcal{H}_K^m$ onto the set $C = \cap_{i=1}^m C_i$, with

$$C_i = \left\{ (\mathbf{z}, f_1, \dots, f_m) : f_i(\mathbf{x}_j) = z_j \ \forall j \in \bar{S}_i, \mathbf{z} \in \mathbb{R}^n, \right. \\ \left. \{f_i\}_{i=1}^m \subset \mathcal{H}_K \right\} \subset \mathbb{R}^n \times \mathcal{H}_K^m. \quad (9)$$

The significance of this observation lies in the fact that the relaxed form of the regularized kernel least-squares estimator has been expressed as a projection onto the intersection of a collection of m convex sets; in particular, note that each set C_i is a subspace. Thus, by Lemma 1, the SOP algorithm can be used to solve the relaxed problem (7). Moreover, computing $P_{C_i}(\cdot)$ requires agent i to gather examples only within its locally accessible database. More precisely, note that for any $\mathbf{v} = (\mathbf{z}, f_1, \dots, f_m) \in \mathbb{R}^n \times \mathcal{H}_K^m$, $P_{C_i}(\mathbf{v}) = (\mathbf{z}^*, f_1^*, \dots, f_m^*)$ where

$$\begin{aligned} f_j^* &= f_j \quad \forall j \neq i \\ f_i^* &= \arg \min_{f \in \mathcal{H}_K} \sum_{j \in \bar{S}_i} (f(x_j) - z_j)^2 + \lambda_i \|f - f_i\|_{\mathcal{H}_K}^2 \\ z_j^* &= z_j \quad \forall j \text{ s.t. } j \notin \bar{S}_i \\ z_j^* &= f_i^*(x_j) \quad \forall j \text{ s.t. } j \in \bar{S}_i \end{aligned}$$

To emphasize, computing $P_{C_i}(\mathbf{v})$ leaves z_j unchanged for all $j \notin \bar{S}_i$ and leaves f_j unchanged for all $j \neq i$. The function associated with agent i , f_i^* can be computed using f_i and S_i after the training data labels $\{y_j\}_{j \in \bar{S}_i}$ have been updated with the corresponding “message variables” $\{z_j\}_{j \in \bar{S}_i}$. Tying these observations together, we are left with an algorithm for collaborative regression estimation which solves a relaxed form of the regularized least-squares estimator (7).

The algorithm is summarized in psuedo-code in Table 1 and depicted pictorially in Figure 5. In words, the algorithm iterates over each agent in turn, allowing them to compute a local kernel estimate and to update the labels in the training database accordingly. Multiple passes (in fact, T cycles) over the agents are made.

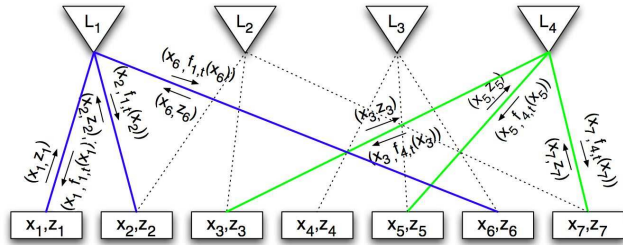


Fig. 5. A Collaborative Training Algorithm

C. Convergence

Note that the asymptotic behavior of the collaborative training algorithm is implied by the analysis of the SOP algorithm. In particular, we have the following.

Theorem 3: Let $(\mathbf{z}, f_{\lambda_1}, \dots, f_{\lambda_m}) \in \mathbb{R}^n \times \mathcal{H}_K^m$ be the solution to (7) and let $\{f_{i,T}\}_{i=1}^m \subset \mathcal{H}_K$ be as defined in the algorithm described in Table I. Then,

$$\lim_{T \rightarrow \infty} f_{i,T} = f_{\lambda_i}$$

for all $i \in \{1, \dots, m\}$.

This theorem follows from Theorem 2 and the fact that convergence in norm implies point-wise convergence in RKHSs.

Given the structure of RKHS and the general analysis in [2], the algorithm is expected to converge linearly for many kernels. We forego a discussion of this important, but technical point for the sake of space.

Observe that Theorem 3 characterizes the output of collaborative training algorithm relative to (7). This characterization is useful insofar as it sheds light on the relationship between the algorithm’s output and (1), the centralized regularized least-squares estimator. The following straightforward generalization of Theorem 1 is a step toward further understanding this important relationship.

Theorem 4: Let $(\mathbf{z}, f_{\lambda_1}, \dots, f_{\lambda_m}) \in \mathbb{R}^n \times \mathcal{H}_K^m$ be the solution to (7). Then, for every agent $i \in \{1, \dots, m\}$, there exists $c_{\lambda_i} \in \mathbb{R}^{|\bar{S}_i|}$ such that

$$f_{\lambda_i}(\cdot) = \sum_{j \in \bar{S}_i} c_{\lambda_{ij}} K(\cdot, x_j). \quad (10)$$

The proof of this theorem follows from the original Representer Theorem (applied to the update equation for $f_{i,t}$) and the fact that \mathcal{H}_K is closed.

The significance of Theorem 4 lies in the fact that the size of any agent’s locally accessible database fundamentally limits the accuracy of that agent’s estimate. In particular, an agent having access to only a few exemplars in an otherwise large training database will still be limited to estimates that lie in the span of functions determined by its local data; thus, local connectivity influences the agent’s bias. Intuitively, however, the message-passing through the training database may optimize the estimator within that limited span if the ensemble is “connected” in some meaningful way. To bear out this intuition in a simplified theoretical setting, we consider a simple notion of connectedness in the next section.

D. A Simplified Setting

For a given ensemble, kernel pair $(\{S_i\}_{i=1}^m, K)$, let us construct an auxiliary graph as follows: let there be a node for every learning agent and let there be an edge between node (i.e., agent) i and node k if the following condition holds:

$$\begin{aligned} \text{span}(\{K(\cdot, x_j)\}_{j \in \bar{S}_i}) &= \text{span}(\{K(\cdot, x_j)\}_{j \in \bar{S}_k}) \\ &= \text{span}(\{K(\cdot, x_j)\}_{j \in \bar{S}_i \cap \bar{S}_k}) \end{aligned} \quad (11)$$

In other words, an edge connects two nodes if the training examples they share determine the space of functions their estimates lie in as dictated by Theorem 4.

Definition 1: Let us call the ensemble, kernel pair $(\{S_i\}_{i=1}^m, K)$ *connected* if and only if the auxiliary graph so constructed is connected.

This definition leads to the following theorem, which can be viewed as a straightforward generalization of Lemma 1.

Theorem 5: Let $(\{S_i\}_{i=1}^m, K)$ be connected and suppose the ensemble employs the collaborative training algorithm using $\{\lambda_i\}_{i=1}^m$. Finally, let f_λ denote the solution to (1) for $\lambda = \sum_{i=1}^m \lambda_i$. Then,

$$f_\lambda = \lim_{T \rightarrow \infty} f_{i,T} \quad (12)$$

for all $i \in \{1, \dots, m\}$.

Theorem 5 follows from Theorem 3 after noting that connectedness implies that the solution to (7) $(\mathbf{z}, f_{\lambda_1}, \dots, f_{\lambda_m})$ satisfies $f_{\lambda_1} = \dots = f_{\lambda_m}$. To illustrate the significance of Theorem 5

Init:	Agents agree on a positive semi-definite kernel $K(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. Training database $S = \{(x_i, z_i)\}_{i=1}^n$ is initialized so that $z_i = y_i \forall i \in \{1, \dots, n\}$.
Train:	for $t = 1, \dots, T$ for $i = 1, \dots, m$ Agent i : Retrieves database $S_i \subseteq S$ Computes $f_{i,t} := \arg \min_{f \in \mathcal{H}_K} \left[\sum_{j \in S_i} (f(\mathbf{x}_j) - z_j)^2 + \lambda_i \ f - f_{i,t-1}\ _{\mathcal{H}_K}^2 \right]$ Updates database: $z_j \leftarrow f_{i,t}(\mathbf{x}_j) \forall (x_j, z_j) \in S_i$ end end end

TABLE I
AN ALGORITHM FOR TRAINING COLLABORATIVELY

and to tie it to the foregoing discussion, consider the following example.

Example 1: Suppose $\mathcal{X} = \mathbb{R}^d$ and that $K(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$ is the linear kernel; in this case, \mathcal{H}_K is the set of linear functions on \mathcal{X} . If $\{S_i\}_{i=1}^m$ is an ensemble with public database of d linearly independent examples (depicted in Figure 3 and discussed in Section III), then $(\{S_i\}_{i=1}^m, K)$ is connected. Therefore, by Theorem 5, the collaborative training algorithm would allow agent i to find the best linear fit to the entire data set S (for the particular choice of regularization parameter λ), despite the fact that only $\frac{n_i+d}{\sum_{i=1}^m n_i+d}$ percent of the data is locally accessible. More generally, if a p^{th} order polynomial kernel is used, then an analogous observation holds when d^p examples are shared.

In this simple example, the potential utility of the collaborative training algorithm is revealed. Consider the extreme case when each agent has access to only a single example in addition to the public database. As the number of agents $m \rightarrow \infty$, the collaborative training algorithm would allow every agent a consistent estimate of the optimal linear least-squares estimate as long as $\sum_{i=1}^m \lambda_i \rightarrow 0$; this is true despite the fact that each agent retains local access to only $d+1$ examples for all m .

IV. DISCUSSION

As described in Table 1, the inner loop of the collaborative training algorithm iterates over agents in the ensemble serially. Note that the ordering is non-essential and parallelism may be introduced. In fact, two agents can train simultaneously as long as they do not share exemplars in their locally accessible training database. In practical settings, multiple-access algorithms that are frequently studied in the communications literature (e.g., ALOHA) may be adapted to negotiate an ordering in a distributed fashion. Since the SOP algorithm and Theorem 2 have been generalized to a very general class of (perhaps random) control orderings [2], Theorem 3 can be extended in many cases. Experiments that validate the collaborative training algorithm in a WSN setting can be found in [14].

In this paper, we have focused exclusively on regularized kernel least-squares regression. However using Bregman’s algorithm [4], the method and many of the theorems may be extended to more general loss functions and regularizers including Bregman divergences.

Those familiar with LDPC codes or Bayes networks may find the current model and algorithm reminiscent of message-passing algorithms such as belief-propagation which are frequently studied in those fields; variational interpretations of kernel methods in the context of Gaussian processes further suggests a relationship

between these works. Formalizing such a connection would likely require one to interpret our “relaxation” in the context of dependency structures in Gaussian processes, and to connect alternating projection algorithms with the generalized distributive law [1].

REFERENCES

- [1] S. M. Aji and R. J. McEliece, “The generalized distributive law,” *IEEE Transactions on Information Theory*, vol. 46, no. 2, pp. 325–343, March 2000.
- [2] H. H. Bauschke and J. M. Borwein, “On projection algorithms for solving convex feasibility problems,” *SIAM Review*, vol. 38, no. 3, pp. 367–426, September 1996.
- [3] A. Bordes, S. Ertekin, J. Weston, and L. Bottou, “Fast kernel classifiers with online and active learning,” *Journal of Machine Learning Research*, vol. 6, pp. 1579–1619, 2005.
- [4] Y. Censor and S. A. Zenios, *Parallel Optimization: Theory, Algorithms, and Applications*. New York: Oxford, 1997.
- [5] 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.
- [6] S. Gambs, B. Kégl, and E. Aïmeur, “Privacy-preserving boosting,” submitted to *Data Mining and Knowledge Discovery*, 2005.
- [7] C. Guestrin, P. Bodi, R. Thibau, M. Paskin, and S. Madde, “Distributed regression: an efficient framework for modeling sensor network data,” in *IPSN’04: Proceedings of the Third International Symposium on Information Processing in Sensor Networks*. New York: ACM Press, 2004.
- [8] L. Györfi, M. Kohler, A. Krzyzak, and H. Walk, *A Distribution-Free Theory of Nonparametric Regression*. New York: Springer, 2002.
- [9] M. Kearns and H. S. Seung, “Learning from a population of hypotheses,” *Machine Learning*, vol. 18, pp. 255–276, 1995.
- [10] G. Kimeldorf and G. Wahba, “Some results on Tchebycheffian spline functions,” *Journal of Mathematical Analysis and Applications*, vol. 33, pp. 82–95, 1971.
- [11] J. Kivinen and M. K. Warmuth, “Additive versus exponentiated gradient updates for linear prediction,” *Information and Computation*, vol. 132, no. 1, pp. 1–64, 1997.
- [12] 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.
- [13] 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.
- [14] J. B. Predd, S. R. Kulkarni, and H. V. Poor, “Regression in sensor networks: Training distributively with alternating projections,” in *Proceedings of the SPIE Conference on Advanced Signal Processing Algorithms, Architectures, and Implementations XV* (invited), San Diego, CA, July 31 – August 4 2005.
- [15] —, “Consistency in models for distributed learning under communication constraints,” *IEEE Transactions on Information Theory*, vol. 52, no. 1, pp. 52–63, Jan. 2006.
- [16] B. Schölkopf and A. Smola, *Learning with Kernels*, 1st ed. Cambridge, MA: MIT Press, 2002.
- [17] P. K. Varshney, *Distributed Detection and Data Fusion*. New York: Springer, 1996.
- [18] J. von Neumann, *Function Operators II*. Princeton, NJ: Princeton University, 1950.
- [19] G. Wahba, *Spline Models for Observational Data*. Philadelphia: SIAM, 1990.