

Multi-Agent Distributed Lifelong Learning for Collective Knowledge Acquisition

Mohammad Rostami
University of Pennsylvania
Philadelphia, PA, USA
mrostami@seas.upenn.edu

Kyungnam Kim
HRL Laboratories, LLC
Malibu, CA, USA
kkim@hrl.com

Soheil Kolouri
HRL Laboratories, LLC
Malibu, CA, USA
skolouri@hrl.com

Eric Eaton
University of Pennsylvania
Philadelphia, PA, USA
eeaton@cis.upenn.edu

ABSTRACT

Lifelong machine learning methods acquire knowledge over a series of consecutive tasks, continually building upon their experience. Current lifelong learning algorithms rely upon a single learning agent that has centralized access to all data. In this paper, we extend the idea of lifelong learning from a single agent to a network of multiple agents that collectively learn a series of tasks. Each agent faces some (potentially unique) set of tasks; the key idea is that knowledge learned from these tasks may benefit other agents trying to learn different (but related) tasks. Our Collective Lifelong Learning Algorithm (CoLLA) provides an efficient way for a network of agents to share their learned knowledge in a distributed and decentralized manner, while preserving the privacy of the locally observed data. Note that a decentralized scheme is a subclass of distributed algorithms where a central server does not exist and in addition to data, computations are also distributed among the agents. We provide theoretical guarantees for robust performance of the algorithm and empirically demonstrate that CoLLA outperforms existing approaches for distributed multi-task learning on a variety of data sets.

KEYWORDS

Lifelong machine learning; multi-agent collective learning; distributed optimization

1 INTRODUCTION

Collective knowledge acquisition is common throughout different societies, from the collaborative advancement of human knowledge to the emergent behavior of ant colonies [14]. It is the product of individual agents, each with their own interests and constraints, sharing and accumulating learned knowledge over time in uncertain and dangerous real-world environments. Our work explores this scenario within machine learning and in particular, considers learning in a network of lifelong machine learning agents [28].

Recent work in lifelong machine learning [26] has explored the notion of a single agent accumulating knowledge over its lifetime.

Proc. of the 17th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2018), M. Dastani, G. Sukthankar, E. Andre, S. Koenig (eds.), July 2018, Stockholm, Sweden

© 2018 International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). All rights reserved.
<https://doi.org/doi>

Such an individual lifelong learning agent reuses knowledge from previous tasks to improve its learning on new tasks, accumulating an internal repository of knowledge over time. This lifelong learning process improves performance over all tasks, and permits the design of adaptive agents that are capable of learning in dynamic environments. Although current work in lifelong learning focuses on a single learning agent that incrementally perceives all task data, many real-world applications involve scenarios in which multiple agents must collectively learn a series of tasks that are distributed amongst them. Consider the following cases:

- Multi-modal task data could only be partially accessible by each learning agent. For example, financial decision support agents may have access only to a single data view of tasks or a portion of the non-stationary data distribution [11].
- Local data processing can be inevitable in some applications, such as when health care regulations prevent personal medical data from being shared between learning systems [38].
- Data communication may be costly or time consuming. For instance, home service robots must process perceptions locally due to the volume of perceptual data, or wearable devices may have limited communication bandwidth [13].
- As a result of data size or geographical distribution of data centers, parallel processing can be essential. Modern big data systems often necessitates parallel processing in the cloud across multiple virtual agents, i.e. CPUs or GPUs [39].

Inspired by the above scenarios, this paper explores the idea of *multi-agent lifelong learning*. We consider multiple collaborating lifelong learning agents, each facing their own series of tasks, that transfer knowledge to collectively improve task performance and increase learning speed. Note that this paper does not address the privacy considerations that may arise from transferring knowledge between agents. Also, despite potential extendability to parallel processing systems, our focus here is more on collaborative agents that receive sequential tasks. Existing methods in the literature have mostly investigated special cases of this setting for distributed multi-task learning (MTL) [7, 13, 25].

To develop multi-agent distributed lifelong learning, we follow a parametric approach and formulate the learning problem as an on-line MTL optimization over a network of agents. Each agent seeks to learn parametric models for its own series of (potentially unique) tasks. The network topology imposes communication constraints

among the agents. For each agent, the corresponding task model parameters are represented as a task-specific sparse combination of atoms of its local knowledge base [15, 22, 26]. The local knowledge bases allow for knowledge transfer from learned tasks to the future tasks for each individual agent. The agents share their knowledge bases with their neighbors, update them to incorporate the learned knowledge representations of their neighboring agents, and come to a local consensus to improve learning quality and speed. Our contribution is to use the Alternating Direction Method of Multipliers (ADMM) algorithm [4] to solve this global optimization problem in an online distributed setting; our approach decouples this problem into local optimization problems that are individually solved by the agents. ADMM allows for transferring the learned local knowledge bases without sharing the specific learned model parameters among neighboring agents. We propose an algorithm with nested loops to allow for keeping the procedure both online and distributed. We call our approach the Collective Lifelong Learning Algorithm (CoLLA). We provide theoretical analysis on the convergence of CoLLA and empirically validate the practicality of the proposed algorithm on variety of standard MTL benchmark datasets.

2 RELATED WORK

This paper considers scenarios where multiple lifelong learning agents learn a series of tasks distributed among them. Each agent shares high level information with its neighboring agents, while processing data privately. Our approach draws upon various sub-fields of machine learning, which we briefly survey below.

Multi-Task and Lifelong Learning: Multi-task learning (MTL) [5] seeks to share knowledge among multiple related tasks. Compared to single task learning (STL), MTL increases generalization performance and reduces the data requirements for learning tasks through benefiting from task similarities. One major challenge in MTL is modeling task similarities to selectively share information and enable knowledge transfer between tasks [5]. If this process identifies incorrect task relationships, sharing knowledge can degrade performance through the phenomenon of negative transfer. Various techniques have been developed to model task relations, including modeling a task distance metric [3], using correlations to determine when transfer is appropriate [33], and regularizing task parameters [1]. An effective parametric approach is to group similar tasks by assuming that task parameters can be represented sparsely in a dictionary domain. Then, by imposing sparsity on task-specific parameters, similar tasks can be grouped together enabling positive knowledge transfer and the learned dictionary would model task relations [15]. Upon learning the dictionary, similar tasks would share a subset of dictionary columns which helps to avoid negative transfer.

Lifelong learning is a special case of MTL in which an agent receives tasks consecutively and continuously and learns the tasks in an online sequential setting. To improve learning performance on each new task, the agent transfers knowledge obtained from the previous tasks, and then stores new or revised knowledge for future use [24]. As a result, the learning agent is able to adapt itself to dynamic environments and drifts in data distribution, and consistently improves its performance. This ability is essential for emerging applications such as personal intelligent robot assistants

and chatbots which continuously interact with many different users. Pioneer works in lifelong learning are built upon extending MTL schemes to an online setting. Ruvolo and Eaton [26] extended the MTL method proposed by Kumar and Daume III [15] to a lifelong learning setting, creating an efficient and fast algorithm for lifelong learning. Our approach is partially based upon their formulation, which serves as the foundation to develop our novel collective lifelong learning framework. Note that unlike our work, most prior MTL and lifelong learning work consider the case where all tasks are accessible by a single agent in a centralized scheme.

Distributed Machine Learning: There has been a growing interest in developing scalable learning algorithms using distributed optimization [40], motivated by the emergence of big data [6], security and privacy constraints [37], and the notion of cooperative and collaborative learning agents [8]. Distributed machine learning allows multiple agents to collaboratively mine information from large-scale data. The majority of these settings are graph-based, where each node in the graph represents a portion of data or an agent. Communication channels between the agents then can be modeled via edges in the graph. Some approaches assume there is a central server (or a group of server nodes) in the network, and the worker agents transmit locally learned information to the server(s), which then perform knowledge fusion [35]. Other approaches assume that processing power is distributed among the agents, which exchange information with their neighbors during the learning process [7]. We formulate our problem in the latter setting, as it is less restrictive. Following the dominant paradigm of distributed optimization, we also assume that the agents are synchronous.

These algorithms formulate learning as an optimization problem over the network and use techniques from distributed optimization to acquire the global solution. Various techniques have been explored, including stochastic gradient descent [35], proximal gradients [17], and ADMM [35]. Within the ADMM framework, it is assumed that the objective function over the network can be decoupled into a sum of independent local functions for each node (usually risk functions) [20], constrained by network topology. Through a number of iterations on primal and dual variables of the Lagrangian function, each node solves a local optimization, and then through information exchange, constraints imposed by the network are realized by updating the dual variable. In scenarios where maximizing a cost for some agents translate to minimizing the cost for others (e.g., adversarial games), game theoretical notions are used to define a global optimal state for the agents [18].

Distributed Multi-task Learning: Although it seems natural to consider MTL agents that collaborate on related tasks, most prior distributed learning work focuses on the setting where all agents try to learn a single task. Only recently have MTL scenarios been investigated where the tasks are distributed [2, 13, 19, 21, 32, 34]. In such a setting, data must not be transferred to a central node because of communication and privacy/security constraints. Only the learned models or high-level information can be exchanged by neighboring agents. Distributed MTL has also been explored in reinforcement learning settings [9] where the focus is on developing a scalable multitask policy search algorithm. These distributed MTL methods are mostly limited to off-line (batch) settings where each agent handles only one task [21, 32]. Jin et al. [13] consider an online setting, but require the existence of a central node, which is

restrictive. In contrast, our work considers decentralized and distributed multi-agent MTL in a lifelong learning setting, without the need for a central server. Moreover, our approach employs homogeneous agents that collaborate to improve collective performance over consecutive distributed tasks in a lifelong learning setting. This can be considered as a special case of concurrent learning where learning task concurrently by multiple agents can speed up learning rate [12].

Similar to prior works [9, 21, 32], we use distributed optimization to tackle the collective lifelong learning problem. These existing approaches can only handle an off-line setting where all the task data is available in a batch for each agent. In contrast, we propose an online learning procedure which can address consecutively arriving tasks. In each iteration, the agents receive and learn their local task models. Since the agents are synchronous, once the tasks are learned by agents, a message-passing scheme is then used to transfer and update knowledge between the neighboring agents in each iteration. In this manner, knowledge will disseminate among all agents over time, improving collective performance. Similar to most distributed learning settings, we assume there is a latent knowledge base that underlies all tasks and that each agent is trying to learn a local version of that knowledge base based on its own (local) observations and knowledge exchange with neighboring agents, modeled by edges (links) of the representing network graph.

3 LIFELONG MACHINE LEARNING

We consider a set of T related (but different) supervised regression or classification tasks, each with labeled training data, i.e. $\{\mathcal{Z}^{(t)} = (X^{(t)}, \mathbf{y}^{(t)})\}_{t=1}^T$, where $X^{(t)} = [\mathbf{x}_1, \dots, \mathbf{x}_M] \in \mathbb{R}^{d \times M}$ represents M task data instances characterized by d features, and $\mathbf{y}^{(t)} = [y_1, \dots, y_m]^T \in \mathcal{Y}^M$ are the corresponding targets. Typically, $\mathcal{Y} = \{\pm 1\}$ for binary classification tasks and $\mathcal{Y} = \mathbb{R}$ for regression tasks. We assume that for each task, the hidden mapping from each data point \mathbf{x}_m to the corresponding target y_m , $f: \mathbb{R}^d \rightarrow \mathcal{Y}$ can be parametrized as $y_m = f(\mathbf{x}_m; \theta^{(t)})$, where $\theta^{(t)} \in \mathbb{R}^d$. In this work, we consider a linear mapping $f(\mathbf{x}_m; \theta^{(t)}) = \langle \theta^{(t)}, \mathbf{x}_m \rangle$ where $\theta^{(t)} \in \mathbb{R}^d$, but our framework is readily generalizable to nonlinear parametric mappings (e.g., using generalized dictionaries [31]). After receiving a task $\mathcal{Z}^{(t)}$, the goal of the agent is to learn the mapping $f(\mathbf{x}_m; \theta^{(t)})$ by estimating the corresponding optimal task parameter $\theta^{(t)}$ using the training data such that it well-generalizes on testing data points from that task. An agent can learn the task models by solving for the optimal parameters $\Theta^* = [\theta^{(1)}, \dots, \theta^{(T)}]$ in the following expected risk minimization (ERM) problem:

$$\min_{\Theta} \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X^{(t)} \sim \mathcal{D}^{(t)}} \left(\mathcal{L} \left(X^{(t)}, \mathbf{y}^{(t)}; \theta^{(t)} \right) \right) + \Omega(\Theta), \quad (1)$$

where $\mathcal{L}(\cdot)$ is a loss function for measuring data fidelity, $\mathbb{E}(\cdot)$ denotes the expectation on the task's data distribution $\mathcal{D}^{(t)}$, and $\Omega(\cdot)$ is a regularization function that models task relations by coupling model parameters to transfer knowledge among the tasks. Almost all parametric MTL, online, and lifelong learning algorithms solve instances of (1) given a particular form of $\Omega(\cdot)$ to impose a specific coupling scheme and an optimization mode, i.e. online or offline.

To model task relations, the GO-MTL algorithm [15] uses classic Empirical Risk Minimization (ERM) to estimate the expected loss and solves the objective (1) by assuming that the task parameters can be decomposed into a shared knowledge dictionary base $L \in \mathbb{R}^{d \times u}$ to facilitate knowledge transfer and task-specific sparse coefficients $\mathbf{s}^{(t)} \in \mathbb{R}^u$, such that $\theta^{(t)} = L\mathbf{s}^{(t)}$. In this factorization, the hidden structure of the tasks is represented in the knowledge dictionary base and similar tasks are grouped by imposing sparsity on $\mathbf{s}^{(t)}$'s. Tasks that use the same columns of the dictionary are clustered to be similar, while tasks that do not share any column can be considered belonging to different groups. In other words, more overlap on sparsity patterns of two tasks imply more similarity between those two tasks. This factorization has been analytically shown to enable knowledge transfer when dealing with related tasks by grouping the similar tasks [15, 22]. Following this assumption and employing ERM, the objective (1) can be estimated as:

$$\min_{L, S} \frac{1}{T} \sum_{t=1}^T \left[\hat{\mathcal{L}} \left(X^{(t)}, \mathbf{y}^{(t)}, L\mathbf{s}^{(t)} \right) + \mu \|\mathbf{s}^{(t)}\|_1 \right] + \lambda \|L\|_F^2, \quad (2)$$

where $S = [\mathbf{s}^{(1)} \dots \mathbf{s}^{(T)}]$ is the matrix of sparse vectors, $\hat{\mathcal{L}}(\cdot)$ is the empirical loss function on task training data, $\|\cdot\|_F$ is the Frobenius norm to regularize complexity and impose uniqueness, $\|\cdot\|_1$ denotes the L_1 norm to impose sparsity on $\mathbf{s}^{(t)}$, and μ and λ are regularization parameters. Eq. (2) is not a convex problem in its general form, but for a convex loss function is convex in each individual optimization variables L and S . Given all tasks' data in a batch, Eq. (2) can be solved in an offline batch-mode by an iterative alternating optimization scheme [15]. In each alternation step, Eq. (2) is solved to update a single variable by treating the other variable to be constant. This scheme leads to an MTL algorithm that enables the agent to share information selectively among the tasks.

Solving Eq. (2) in an offline setting is not applicable for lifelong learning because all tasks are not available in a single batch. A lifelong learning agent faces tasks sequentially over its lifetime, where each task should be learned efficiently and fast using knowledge transferred from past experience. The agent does not know the total number of tasks, nor the order the tasks are received. In other words, for each task $\mathcal{Z}^{(t)}$, the corresponding parameter $\theta^{(t)}$ is learned using knowledge obtained from the training data from previous tasks $\{\mathcal{Z}^{(1)}, \dots, \mathcal{Z}^{(t-1)}\}$ and upon learning $\mathcal{Z}^{(t)}$, the learned or updated knowledge is stored to benefit learning future tasks in an online scheme. To solve Eq. (2) in an online lifelong learning setting, Ruvolo and Eaton [26] first approximate the loss function $\mathcal{L}(X^{(t)}, \mathbf{y}^{(t)}, L\mathbf{s}^{(t)})$ using a second-order Taylor expansion of the loss function around a single-task ridge-optimal parameters. This technique reduces the objective function of Eq. (2) to the problem of online dictionary learning [20] as follows:

$$\min_L \frac{1}{T} \sum_{t=1}^T F^{(t)}(L) + \lambda \|L\|_F^2, \quad (3)$$

$$F^{(t)}(L) = \min_{\mathbf{s}^{(t)}} \left[\|\alpha^{(t)} - L\mathbf{s}^{(t)}\|_{\Gamma^{(t)}}^2 + \mu \|\mathbf{s}^{(t)}\|_1 \right], \quad (4)$$

where $\|\mathbf{x}\|_A^2 = \mathbf{x}^T A \mathbf{x}$, $\alpha^{(t)} \in \mathbb{R}^d$ is the ridge estimator for task $\mathcal{Z}^{(t)}$:

$$\alpha^{(t)} = \arg \min_{\theta^{(t)}} \left[\hat{\mathcal{L}}(\theta^{(t)}) + \gamma \|\theta^{(t)}\|_2^2 \right], \quad (5)$$

and $\Gamma^{(t)}$ is the Hessian of the loss $\hat{\mathcal{L}}(\cdot)$ at $\alpha^{(t)}$ which is assumed to be strictly positive definite. Also, $\gamma \in \mathbb{R}^d$ is the ridge regularization parameter. To solve Eq. (3) in an online setting, still an alternation scheme is used but when a new task arrives, only the corresponding sparse vector $s^{(t)}$ for that task is computed using L to update the sum $\sum_{t=1}^T F(L)$. In this setting, Eq. (4) becomes a task-specific online operation that leverages knowledge transfer. Finally the shared basis L is updated via Eq. (3) to store the learned knowledge from $\mathcal{Z}^{(t)}$ for future use. Despite using Eq. (4) as an approximation to solve for s , Ruvolo and Eaton [26] proved that the learned knowledge base L stabilize as more tasks are learned and would eventually converge to the offline scheme solution of Kumar and Daume III [15]. Moreover, the solution of Eq. (1) converges almost surely to the solution of Eq. (2) as $T \rightarrow \infty$. While this technique leads to an efficient and fast algorithm for lifelong learning, it requires centralized access to all tasks' data by a single agent. The approach we explore, COLLA, benefits from the idea of the second-order Taylor approximation and online optimization scheme proposed by Ruvolo and Eaton [26], but eliminates the need for a centralized data access. CoLLA achieves a distributed and decentralized knowledge update by formulating a multi-agent lifelong learning optimization problem over a network of collaborating lifelong learning agents. The resulting optimization problem, can be solved in a distributed setting, enabling collective learning, as we describe next.

4 MULTI-AGENT LIFELONG LEARNING

We consider a network of N collaborating lifelong learning agents with an arbitrary order on the agents. Note however this arbitrary order is restrictive and needs to be known and fixed. Each agent receives a potentially unique tasks at each time-step in a lifelong learning setting. There is also some true underlying hidden knowledge base for all tasks, and each agent learns a local view of this knowledge base based on its own task distribution. We assume that each agent i solves a local version of the objective (3) to estimate its own local knowledge base L_i . We also assume that the agents are synchronous, meaning that at each time step, they simultaneously receive and learn one task. We model the communication mode of these agents by an undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where the set of static nodes $\mathcal{V} = \{1, \dots, N\}$ denotes the agents and the set of edges $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$, with $|\mathcal{E}| = e$, specifies possibility of communication between the agents. We assume for each edge $(i, j) \in \mathcal{E}$, the nodes i and j are connected or they can communicate information, with $j > i$ for uniqueness and set orderability. We define the neighborhood $\mathcal{N}(i)$ of node i as the set of all nodes that are connected to it. To allow for knowledge and information flow between all the agents, we further assume that the network graph is connected. Note that all nodes are assumed to be homogeneous and as a result there is no central node that guides collaboration among the agents.

We use the graph structure to formulate a lifelong machine learning problem on this network. Although each agent learns its own individual dictionary, we encourage local dictionaries of neighboring nodes (agents) to be similar by adding a set of soft equality constraints on neighboring dictionaries, i.e. $L_i = L_j, \forall (i, j) \in \mathcal{E}$. We can represent all these constraints as a single linear combination on local dictionaries. It is easy to show these e equality constraints can be written compactly as $(H \otimes I_{d \times d})\tilde{L} = \mathbf{0}_{ed \times u}$, where $H \in \mathbb{R}^{e \times N}$

is the node arc-incident matrix of \mathcal{G} (for a given row $1 \leq l \leq e$, corresponding to the l^{th} edge (i, j) , $H_{lq} = 0$ except for $H_{li} = 1$ and $H_{lj} = -1$), $I_{d \times d}$ is the identity matrix, $\tilde{L} = [L_1^\top, \dots, L_N^\top]^\top$, and \otimes denotes the Kronecker product. Let $E_i \in \mathbb{R}^{ed \times d}$ be a column partition of $E = (H \otimes I_d) = [E_1, \dots, E_N]$; we can compactly write the e equality constraints as $\sum_i E_i L_i = \mathbf{0}_{ed \times u}$. Each of $E_i \in \mathbb{R}^{ed \times d}$ matrices is a tall block matrix consisting of $d \times d$ blocks, $\{[E_i]_j\}_{j=1}^e$, that are either zero matrix $\forall j \notin \mathcal{N}(i)$, $I_d, \forall j \in \mathcal{N}(i), j > i$, or $-I_d, \forall j \in \mathcal{N}(i), j < i$. Note that $E_i^\top E_j = O_d$ if $j \notin \mathcal{N}(i)$, where is $d \times d$ zero matrix. Following this notation, we can reformulate the MTL objective (3) for multiple agents as the following linearly constrained optimization problem over the network graph \mathcal{G} :

$$\begin{aligned} \min_{L_1, \dots, L_N} \quad & \frac{1}{T} \sum_{t=1}^T \sum_{i=1}^N F_i^{(t)}(L_i) + \lambda \|L_i\|_F^2 \\ \text{s.t.} \quad & \sum_{i=1}^N E_i L_i = \mathbf{0}_{ed \times u} \end{aligned} \quad (6)$$

Note that in Eq. (6) optimization variables are not coupled by a global variable and hence in addition to being a distributed problem, Eq. (6) is also a decentralized problem. In order to deal with the dynamic nature and time-dependency of the objective (6), we assume that at each time step t , each agent receives a task and computes $F_i^{(t)}(L_i)$ locally via (4) based on this local task. Then, through K information exchanges during that time step, the local dictionaries are updated such that the agents reach a local consensus and hence benefit from all the tasks that are received by the network in that time step. To split the constrained objective (6) into a sequence of local unconstrained agent-level problems, we use the extended ADMM algorithm [20, 23]. This algorithm generalizes ADMM [4] to account for linearly constrained convex problems with a sum of N separable objective functions. Similar to ADMM, we first need to form the augmented Lagrangian $\mathcal{J}_T(L_1, \dots, L_N, Z)$ for problem (6) at time t in order to replace the constrained problem by an unconstrained objective function which has an added penalty term:

$$\begin{aligned} \mathcal{J}_T(L_1, \dots, L_N, Z) = & \frac{1}{T} \sum_{t=1}^T \sum_{i=1}^N F_i^{(t)}(L_i) + \\ & \lambda \|L_i\|_F^2 + \langle Z, \sum_{i=1}^N E_i L_i \rangle + \frac{\rho}{2} \left\| \sum_{i=1}^N E_i L_i \right\|_F^2, \end{aligned} \quad (7)$$

where $\langle Z, \sum_{i=1}^N E_i L_i \rangle = \text{tr}(Z^\top \sum_{i=1}^N E_i L_i)$ denotes the matrix trace inner product, $\rho \in \mathbb{R}^+$ is a regularization penalty term parameter for violation of the constraint. Finally, the block matrix $Z = [Z_1^\top, \dots, Z_e^\top]^\top \in \mathbb{R}^{ed \times u}$ is the ADMM dual variable. The extended ADMM algorithm solves Eq.(6) by iteratively updating the dual and

primal variables using the following local split iterations:

$$\begin{aligned} \mathbf{L}_1^{k+1} &= \operatorname{argmin}_{\mathbf{L}_1} \mathcal{J}_T(\mathbf{L}_1, \mathbf{L}_2^k, \dots, \mathbf{L}_N^k, \mathbf{Z}^k), \\ \mathbf{L}_2^{k+1} &= \operatorname{argmin}_{\mathbf{L}_2} \mathcal{J}_T(\mathbf{L}_1^{k+1}, \mathbf{L}_2, \dots, \mathbf{L}_N^k, \mathbf{Z}^k), \\ &\vdots \end{aligned} \quad (8)$$

$$\begin{aligned} \mathbf{L}_N^{k+1} &= \operatorname{argmin}_{\mathbf{L}_N} \mathcal{J}_T(\mathbf{L}_1^{k+1}, \mathbf{L}_2^{k+1}, \dots, \mathbf{L}_N, \mathbf{Z}^k), \\ \mathbf{Z}^{k+1} &= \mathbf{Z}^k + \rho \left(\sum_{i=1}^N \mathbf{E}_i \mathbf{L}_i^{k+1} \right). \end{aligned} \quad (9)$$

The first N problems are primal agent-specific problems to update each local dictionary and the last problem updates the dual variable. These iterations split the objective (7) into local primal optimization problems to update each of the \mathbf{L}_i 's, and then synchronize the agents to share information through updating the dual variable.

Note that the j 'th column of \mathbf{E}_i is only nonzero when $j \in \mathcal{N}(i)$ [$\mathbf{E}_i]_j = \mathbf{0}_d, \forall j \notin \mathcal{N}(i)$], hence the update rule for the dual variable is indeed e local block updates by agents that share an edge:

$$\mathbf{Z}_i^{k+1} = \mathbf{Z}_i^k + \rho(\mathbf{L}_i^{k+1} - \mathbf{L}_i^k), \quad (10)$$

for the l^{th} edge (i,j) . This means that to update the dual variable, agent i solely needs to keep track of copies of those blocks \mathbf{Z}_l that are shared with neighboring agents, reducing (9) to a set of distributed local operations. Note that iterations in (8) and (10) are performed K times at each instance t for each agent to allow for agents to converge to a stable solution. At each time step t , the stable solution from the previous time step $t-1$ is used to initialize dictionaries and the dual variable in (8). Due to convergence guarantees of extended ADMM [20], this simply means that at each iterations all the tasks that are received by the agents are considered to update the knowledge bases.

4.1 Dictionary Update Rule

Splitting an optimization using ADMM is particularly helpful if optimization on primal variables can be solved efficiently, e.g. closed form solution. We show that the local primal updates in (8) can be solved in closed form. We simply compute and then null the gradients of the primal problems, which leads to systems of linear problems for each local dictionary \mathbf{L}_i :

$$\begin{aligned} \frac{\partial \mathcal{J}_T}{\partial \mathbf{L}_i} &= \frac{1}{T} \sum_t \frac{\partial \mathcal{F}^{(t)}}{\partial \mathbf{L}_i} + \\ &\rho \mathbf{E}_i^\top (\mathbf{E}_i \mathbf{L}_i + \sum_{j>i} \mathbf{E}_j \mathbf{L}_j^k + \sum_{j<i} \mathbf{E}_j \mathbf{L}_j^{k+1} + \frac{1}{\rho} \mathbf{Z}) + 2\lambda \mathbf{L}_i \\ &= \frac{2}{T} \sum_{t=1}^T \Gamma_i^{(t)} (\mathbf{L}_i \mathbf{s}_i^{(t)} - \boldsymbol{\alpha}_i^{(t)}) \mathbf{s}_i^{(t)\top} + \\ &\mathbf{E}_i^\top (\mathbf{E}_i \mathbf{L}_i + \sum_{j>i} \mathbf{E}_j \mathbf{L}_j^k + \sum_{j<i} \mathbf{E}_j \mathbf{L}_j^{k+1} + \frac{1}{\rho} \mathbf{Z}) + 2\lambda \mathbf{L}_i \\ &= \frac{2}{T} \sum_t \Gamma_i^{(t)} \mathbf{L}_i \mathbf{s}_i^{(t)} \mathbf{s}_i^{(t)\top} + \\ &(\rho |\mathcal{N}(i)| + 2\lambda) \mathbf{L}_i + \sum_{j \in \mathcal{N}(i)} \mathbf{E}_i^\top \mathbf{Z}_j - \frac{2}{T} \sum_t \Gamma_i^{(t)} \boldsymbol{\alpha}_i^{(t)} \mathbf{s}_i^{(t)\top} \\ &\rho \left(\sum_{j<i, j \in \mathcal{N}(i)} \mathbf{E}_i^\top \mathbf{E}_j \mathbf{L}_j^{k+1} + \sum_{j>i, j \in \mathcal{N}(i)} \mathbf{E}_i^\top \mathbf{E}_j \mathbf{L}_j^k \right) = 0. \end{aligned} \quad (11)$$

Note that despite our compact representation, primal iterations in (8) involve only dictionaries from neighboring agents ($\forall j \notin \mathcal{N}(i)$ because $\mathbf{E}_i \mathbf{E}_j = \mathbf{0}$ and $[\mathbf{E}_i]_j = \mathbf{0}_d, \forall j \notin \mathcal{N}(i)$). Moreover, only blocks of dual variable \mathbf{Z} that correspond to neighboring agents are needed to update each knowledge base. This means that iterations in (11) are also fully distributed and decentralized local operations.

To solve for \mathbf{L}_i , we vectorize both sides of Eq. (11) and then after applying a property of Kronecker ($(\mathbf{B}^\top \otimes \mathbf{A}) \operatorname{vec}(\mathbf{X}) = \operatorname{vec}(\mathbf{A} \mathbf{X} \mathbf{B})$) on line 5 of the equation, Eq. (11) simplifies to the following linear equation update rules for local knowledge base dictionaries:

$$\begin{aligned} \mathbf{A}_i &= \left(\frac{\rho}{2} |\mathcal{N}(i)| + \lambda \right) \mathbf{I}_{dk} + \frac{1}{T} \sum_{t=1}^T (\mathbf{s}_i^{(t)} \mathbf{s}_i^{(t)\top}) \otimes \Gamma_i^{(t)}, \\ \mathbf{b}_i &= \operatorname{vec} \left(\frac{1}{T} \sum_{t=1}^T \mathbf{s}_i^{(t)\top} \otimes (\boldsymbol{\alpha}_i^{(t)\top} \Gamma_i^{(t)}) - \frac{1}{2} \sum_{j \in \mathcal{N}(i)} \mathbf{E}_i^\top \mathbf{Z}_j \right. \\ &\quad \left. - \frac{\rho}{2} \left(\sum_{j<i, j \in \mathcal{N}(i)} \mathbf{E}_i^\top \mathbf{E}_j \mathbf{L}_j^{k+1} + \sum_{j>i, j \in \mathcal{N}(i)} \mathbf{E}_i^\top \mathbf{E}_j \mathbf{L}_j^k \right) \right), \\ \mathbf{L}_i &\leftarrow \operatorname{mat}_{d,k} (\mathbf{A}_i^{-1} \mathbf{b}_i), \end{aligned} \quad (12)$$

where $\operatorname{vec}(\cdot)$ denotes the matrix to vector (via column stacking) and $\operatorname{mat}(\cdot)$ denotes the vector to matrix operations. To avoid the sums on t and storing learned tasks data, we construct both \mathbf{A}_i and \mathbf{b}_i incrementally as tasks are learned. Our method, the Collective Lifelong Learning Algorithm (CoLLA), is summarized in Algorithm 1.

5 THEORETICAL CONVERGENCE GUARANTEES

In this section, we provide a proof for convergence of Algorithm 1. We use techniques from Ruvolo and Eaton [26], adapted originally from Mairal et al. [20] to demonstrate that Algorithm 1 converges to a stationary point of the risk function. In the proof, we assume that the following assumptions hold:

(A) The data distribution has a compact support. This assumption enforces boundedness on $\boldsymbol{\alpha}^{(t)}$ and $\Gamma^{(t)}$, and subsequently on \mathbf{L}_i and $\mathbf{s}^{(t)}$ (see [20] for details).

(B) The LASSO problem in Eq. (4) admits a unique solution according to one of uniqueness conditions for LASSO [29]. As a result, the functions $F_i^{(t)}$ are well-defined.

(C) The matrices $\mathbf{L}_i^\top \Gamma^{(t)} \mathbf{L}_i$ are strictly positive definite. As a result, the functions $F_i^{(t)}$ are all strongly convex.

Our proof involves two steps. First, we show that the inner loop with variable k in Algorithm 1 converges to a consensual solution for all i and all t . Next, we prove that the outer loop on t is also convergent, showing that the collectively learned dictionary stabilizes as more tasks are learned. For the first step, we outline the following theorem on the convergence of extended ADMM algorithm.

THEOREM 5.1. (Theorem 4.1 of Han and Yuan [10])

Suppose we have an optimization problem in the form of Eq. (6), where the functions $g_i(\mathbf{L}_i) := \sum_i F_i^{(t)}(\mathbf{L}_i)$ are strongly convex with modulus η_i . Then, for any $0 < \rho < \min_i \left\{ \frac{2\eta_i}{3(N-1)\|\mathbf{E}_i\|^2} \right\}$, iterations in Eq. (8) and Eq. (9) converge to a solution of Eq. (6)

Algorithm 1 CoLLA (k, d, λ, μ, ρ)

```

1:  $T \leftarrow 0$ ,  $\mathbf{A} \leftarrow \mathbf{zeros}_{kd,kd}$ ,
2:  $\mathbf{b} \leftarrow \mathbf{zeros}_{k,1}$ ,  $\mathbf{L}_i \leftarrow \mathbf{zeros}_{d,k}$ 
3: while MoreTrainingDataAvailable() do
4:    $T \leftarrow T + 1$ 
5:   while  $i \leq N$  do
6:      $(\mathbf{X}_i^{(t)}, \mathbf{y}_i^{(t)}, t) \leftarrow \text{getTrainingData}()$ 
7:      $(\boldsymbol{\alpha}_i^{(t)}, \Gamma_i^{(t)}) \leftarrow \text{singleTaskLearner}(X^t, y^t)$ 
8:      $\mathbf{s}_i^{(t)} \leftarrow \text{Equation 4}$ 
9:     while  $k \leq K$  do
10:       $\mathbf{A}_i \leftarrow \mathbf{A}_i + (\mathbf{s}_i^{(t)} \mathbf{s}_i^{(t)\top}) \otimes \Gamma_i^{(t)}$ 
11:       $\mathbf{b}_i \leftarrow \mathbf{b}_i + \text{vec}(\mathbf{s}_i^{(t)\top} \otimes (\boldsymbol{\alpha}_i^{(t)\top} \Gamma_i^{(t)}))$ 
12:       $\mathbf{L}_i \leftarrow \text{reinitializeAllZero}(\mathbf{L}_i)$ 
13:       $\mathbf{b}_i \leftarrow \frac{1}{T} \mathbf{b}_i + \text{vec}(-\frac{1}{2} \sum_{j \in \mathcal{N}(i)} \mathbf{E}_i^\top \mathbf{Z}_j$ 
14:         $-\frac{\rho}{2} (\sum_{j < i, j \in \mathcal{N}(i)} \mathbf{E}_i^\top \mathbf{E}_j \mathbf{L}_j^{k+1}$ 
15:         $+ \sum_{j > i, j \in \mathcal{N}(i)} \mathbf{E}_i^\top \mathbf{E}_j \mathbf{L}_j^k))$ 
16:       $\mathbf{L}_i^k \leftarrow \text{mat}(\left(\frac{1}{T} \mathbf{A}_i + (\frac{\rho}{2} |\mathcal{N}(i)| + \lambda) \mathbf{I}_{kd}\right)^{-1} \mathbf{b}_i)$ 
17:       $\mathbf{Z}_i^{k+1} = \mathbf{Z}_i^k + \rho(\mathbf{L}_i^{k+1} - \mathbf{L}_i^k)$  (distributed)
18:    end while
19:  end while

```

Note that in Algorithm 1, $F_i^{(t)}(\mathbf{L}_i)$ is a quadratic function of \mathbf{L}_i with a symmetric positive definite Hessian and thus $g_i(\mathbf{L}_i)$ as an average of strongly convex functions is also strongly convex. So the required condition for Theorem 5.1 is satisfied and at each time step, the inner loop on k would converge. We represent the consensual dictionary of the agents after ADMM convergence at time $t = T$ with $\mathbf{L}_T = \mathbf{L}_i|_{t=T}, \forall i$ (solution obtained via Eq. (9) and Eq. (6) at $t = T$) and demonstrate that this matrix becomes stable as t grows (outer loop converges), proving overall convergence of the algorithm. More precisely, \mathbf{L}_T is minimizer of the augmented Lagrangian $\mathcal{J}_T(\mathbf{L}_1, \dots, \mathbf{L}_N, \mathbf{Z})$ at $t = T$ and $\mathbf{L}_1 = \dots = \mathbf{L}_N$. Also note that upon convergence of ADMM, $\sum_i \mathbf{E}_i \mathbf{L}_i = \mathbf{O}$. Hence \mathbf{L}_T is the minimizer of the following risk function, derived from Eq. (7):

$$\hat{\mathcal{R}}_T(\mathbf{L}) = \frac{1}{T} \sum_{t=1}^T \sum_{i=1}^N F_i^{(t)}(\mathbf{L}) + \lambda \|\mathbf{L}\|_F^2. \quad (13)$$

We also use the following lemma in our proof [26]:

LEMMA 5.2. *The function $\hat{\mathcal{Q}}_T(\mathbf{L}) = \hat{\mathcal{R}}_T(\mathbf{L}) - \hat{\mathcal{R}}_{T+1}(\mathbf{L})$ is a Lipschitz function: $\forall \mathbf{L}$ and \mathbf{L}' , $|\hat{\mathcal{Q}}_T(\mathbf{L}') - \hat{\mathcal{Q}}_T(\mathbf{L})| \leq O(\frac{1}{T+1}) \|\mathbf{L}' - \mathbf{L}\|$.*

Proof: Note that after algebraic simplifications, we can conclude: $\hat{\mathcal{Q}}_T(\mathbf{L}) = (\frac{1}{T(T+1)} \sum_{t=1}^T \sum_{i=1}^N F_i^{(t)}(\mathbf{L}) - \frac{1}{T+1} F_i^{(T+1)})$. Now note that the functions $F_i^{(t)}(\mathbf{L})$ are quadratic forms with positive definite Hessian matrices and hence are Lipschitz functions, all with Lipschitz parameters upper-bounded by the largest eigenvalue of all Hessian

matrices. Using the definition for a Lipschitz function, it is easy to demonstrate that $\hat{\mathcal{R}}_T(\cdot)$ is also Lipschitz with Lipschitz parameter $O(\frac{1}{T+1})$, because of averaged quadratic form terms in Eq. (13).

Now we can prove the convergence of Algorithm 1:

LEMMA 5.3. *$\mathbf{L}_{T+1} - \mathbf{L}_T = O(\frac{1}{T+1})$, showing that Algorithm 1 converges to a stable dictionary as T grows large.*

Proof: First note that $\hat{\mathcal{R}}_T(\cdot)$ is a strongly convex function for all T . Let η_T be the strong convexity modulus. From the definition, for two points \mathbf{L}_{T+1} and \mathbf{L}_T , we have: $\hat{\mathcal{R}}_T(\mathbf{L}_{T+1}) \geq \hat{\mathcal{R}}_T(\mathbf{L}_T) + \nabla \hat{\mathcal{R}}_T^\top(\mathbf{L}_T)(\mathbf{L}_{T+1} - \mathbf{L}_T) + \frac{\eta_T}{2} \|\mathbf{L}_{T+1} - \mathbf{L}_T\|_F^2$. Since \mathbf{L}_T is minimizer of $\hat{\mathcal{R}}_T(\cdot)$:

$$\hat{\mathcal{R}}_T(\mathbf{L}_{T+1}) - \hat{\mathcal{R}}_T(\mathbf{L}_T) \geq \frac{\eta_T}{2} \|\mathbf{L}_{T+1} - \mathbf{L}_T\|_F^2. \quad (14)$$

On the other hand, from Lemma 5.2:

$$\begin{aligned} \hat{\mathcal{R}}_T(\mathbf{L}_{T+1}) - \hat{\mathcal{R}}_T(\mathbf{L}_T) &= \hat{\mathcal{R}}_T(\mathbf{L}_{T+1}) - \hat{\mathcal{R}}_{T+1}(\mathbf{L}_{T+1}) + \\ &\hat{\mathcal{R}}_{T+1}(\mathbf{L}_{T+1}) - \hat{\mathcal{R}}_{T+1}(\mathbf{L}_T) + \hat{\mathcal{R}}_{T+1}(\mathbf{L}_T) - \hat{\mathcal{R}}_T(\mathbf{L}_T) \leq \\ &\hat{\mathcal{Q}}_T(\mathbf{L}_{T+1}) - \hat{\mathcal{Q}}_T(\mathbf{L}_T) \leq O(\frac{1}{T+1}) \|\mathbf{L}_{T+1} - \mathbf{L}_T\|. \end{aligned} \quad (15)$$

Note that the first two terms in the second line in the above as a whole is negative since \mathbf{L}_{T+1} is the minimizer of $\hat{\mathcal{R}}_{T+1}$. Now combining (14) and (15), it is easy to show that $\|\mathbf{L}_{T+1} - \mathbf{L}_T\|_F^2 \leq O(\frac{1}{T+1})$, concluding the proof:

$$\|\mathbf{L}_{T+1} - \mathbf{L}_T\|_F^2 \leq O(\frac{1}{T+1}) \quad \blacksquare \quad (16)$$

Thus, Algorithm 1 converges as t increases. We also show that the distance between \mathbf{L}_T and the set of stationary points of the agents' true expected costs $\mathcal{R}_T = \mathbb{E}_{\mathbf{X}^{(t)} \sim \mathcal{D}^{(t)}}(\hat{\mathcal{R}}_T)$ converges almost surely (a.s.) to 0 as $T \rightarrow \infty$. We use two theorems from Mairal et al. [20] for this purpose.

THEOREM 5.4. [20] *Consider the empirical risk function $\hat{q}_T(\mathbf{L}) = \frac{1}{T} \sum_{t=1}^T F^{(t)}(\mathbf{L}) + \lambda \|\mathbf{L}\|_F^2$ with $F^{(t)}$ as defined in (4) and the true risk function $q_T(\mathbf{L}) = \mathbb{E}_{\mathbf{X}^{(t)} \sim \mathcal{D}^{(t)}}(\hat{q}_T(\mathbf{L}))$, and make assumptions (A–C). Then both risk functions converge a.s. as $\lim_{T \rightarrow \infty} \hat{q}_T(\mathbf{L}) - q_T(\mathbf{L}) = 0$.*

Note that we can apply this theorem on \mathcal{R}_T and $\hat{\mathcal{R}}_T$, because the inner sum in (13) does not violate the assumptions of Theorem 5.4. This is because the functions $g_i(\cdot)$ are all well-defined and are strongly convex with strictly positive definite Hessians (the sum of positive definite matrices is positive definite). Thus, $\lim_{T \rightarrow \infty} \hat{\mathcal{R}}_T - \mathcal{R}_T = 0$ a.s.

THEOREM 5.5. [20] *Under assumptions (A–C), the distance between the minimizer of $\hat{q}_T(\mathbf{L})$ and the stationary points of $q_T(\mathbf{L})$ converges almost surely to zero.*

Again, this theorem is applicable on \mathcal{R}_T and $\hat{\mathcal{R}}_T$ and thus Algorithm 1 converges to a stationary point of the true risk.

Computational Complexity At each time-step, each agent computes the optimal ridge parameter $\boldsymbol{\alpha}^{(t)}$ and the Hessian matrix $\Gamma^{(t)}$ for the received task. This has a cost of $O(\xi(d, M))$, where $\xi(\cdot)$ depends on the base learner. The cost of updating \mathbf{L}_i and $\mathbf{s}_i^{(t)}$ alone is $O(u^2 d^3)$ [26], and so the cost of updating all local dictionaries by the agents is $O(Nu^2 d^3)$. Note that this step is performed K times in each time-step. Finally, updating the dual variable requires a cost of eud . This leads to overall cost of $O(N\xi(d, M) + K(Nu^2 d^3 + eud))$, which

is independent of T but accumulates as more tasks are learned. We can think of the factor K in the second term as communication cost because in a centralized scheme we would not need these repetitions which requires sharing the local bases with the neighbors. Also, note that if the number of data points per task is big enough, it certainly is more costly to send data to a single server and learn the tasks in a centralized optimization scheme.

6 EXPERIMENTAL RESULTS

We validate our algorithm by comparing it against: 1) STL, a lower-bound to measure the effect of positive transfer from the learned tasks, 2) ELLA, to demonstrate that collaboration between the agents improves overall performance compared to ELLA, 3) offline CoLLA, as an upper-bound to our online distributed algorithm, and finally 4) GO-MTL, as an absolute upper-bound (since GO-MTL is a batch MTL method) to assess the performance of our algorithm from different perspectives. Throughout all experiments, we present and compare the average performance of all agents.

6.1 Datasets

We used four benchmark MTL datasets in our experiments, including two classification and two regression datasets: 1) land mine detection in radar images [36], 2) facial expression identification from photographs of a subject's face [30], 3) predicting London students' scores using school-specific and student-specific features [1], and 4) predicting ratings of customers for different computer models [16]. Below we describe each dataset. Note that for each dataset, we assume that the tasks are distributed equally among the agents. We used different numbers of agents across the datasets to explore various network sizes of the multi-agent system in our framework.

Land Mine Detection: This dataset consists of binary classification tasks to detect whether a geographical region contains a land mine from a radar image. There are 29 tasks, each corresponding to a different geographical region, with a total 14,820 data points. Each data point consists of nine features, including four moment-based features, three correlation-based, one energy-ratio, and one spatial variance feature (see Xue et. al. [36] for details), all automatically extracted from radar images. We added a bias term as a 10^{th} feature. The dataset has a natural dichotomy between foliated and desert regions. We assumed there are two collaborating agents, each dealing solely with one geographical region type.

Facial Expression Recognition: This dataset consists of binary facial expression recognition tasks [30]. We followed Ruvolo and Eaton [26] and chose tasks detecting three facial action units (upper lid raiser, upper lip raiser, and lip corner pull) for seven different subjects, resulting in 21 total tasks, each with 450–999 data points. A Gabor pyramid scheme is used to extract a total number of 2,880 Gabor features from images of a subject's face (see Ruvolo and Eaton [26] for details). Each data point consists of the first 100 PCA components of these Gabor features. We used three agents, each of which learns seven randomly selected tasks. Given that facial expression recognition is a core task for personal assistant robots, each agent can be considered a personal service robot that interacts with few people in a specific environment.

London Schools: This dataset was provided by the Inner London Education Authority. It consists of examination scores of 15,362

students (each assumed to be a data point) in 139 secondary schools (each assumed to be a single task) during three academic years. The goal is to predict the score of students of each school using provided features as a regression problem. We used the same 27 categorical features as described by Kumar and Daume III [15], consisting of eight school-specific features and 19 student-specific features, all encoded as binary features. We also added a feature to account for the bias term. For this dataset, we considered six agents and allocated 23 tasks randomly to each agent.

Computer Survey: The goal in this dataset is to predict the likelihood of purchasing one of 20 different computers by 190 subjects. Each subject is assumed to be a task and its ratings determines the task data points. Each data point consists of 13 binary features, e.g. guarantee, telephone hot line, etc. (see Lenk et al. [1996] for details). We added a feature to account for the bias term. The output is a rating on a scale 0–10 collected in a survey from the subjects. We considered 19 agents and randomly allocated ten tasks to each.

6.2 Evaluation methodology

For each experiment, we randomly split the data for each task evenly into training and testing sets. We performed 100 learning trials on training sets and reported the average performance on the testing sets for these trials as well as the performance variance. For the online settings (CoLLA and ELLA), we also randomized the order of task presentation in each trial to rule out biases in the order of learning tasks. For the offline settings (GO-MTL, Dist. CoLLA, STL), we reported the average asymptotic performance on all task because all tasks are presented and learned simultaneously. We used brute force search to cross validate the parameters u , λ , μ , and ρ for each dataset; these parameters were selected to maximize the performance on a validation set for each algorithm independently. Parameters λ , μ , and ρ are selected from the set $\{10^n | -6 \leq n \leq 6\}$ and u from $\{1, \dots, \max(10, \frac{T}{4})\}$ (Note that $u \ll T$).

Quality of agreement among the agents: The inner loop in Algorithm 1 implements information exchange between the agents. For effective collective learning, agents need to come to an agreement at each time-step which is guaranteed by ADMM if K is chosen large enough. During our experiments, we noticed that initially K needs to be fairly large but as more tasks are learned, it can be decreased over time $K \propto K_1 + K_2/t$ without considerable change in performance ($K_1 \in \mathbb{N}$ is generally small and $K_2 \in \mathbb{N}$ is large). This is expected because the learned tasks by all agents are homogeneous and hence as more tasks are learned, knowledge transfer from previously learned tasks makes local dictionaries closer.

For the two regression problems, we used standard root mean-squared error (RMSE) on the testing set to measure performance of the algorithms. For the two classification problems, we used the area under the ROC curve (AUC) to measure performance. We used AUC because both classification datasets have highly unbiased distributions, making RMSE less informative. Unlike AUC, RMSE is agnostic to the trade-off between false-positives and false-negatives, which can vary in terms of importance in different applications.

6.3 Results

We clarify that In the result section, the number of learned tasks is equal for both CoLLA and ELLA. While the x-axis is the number

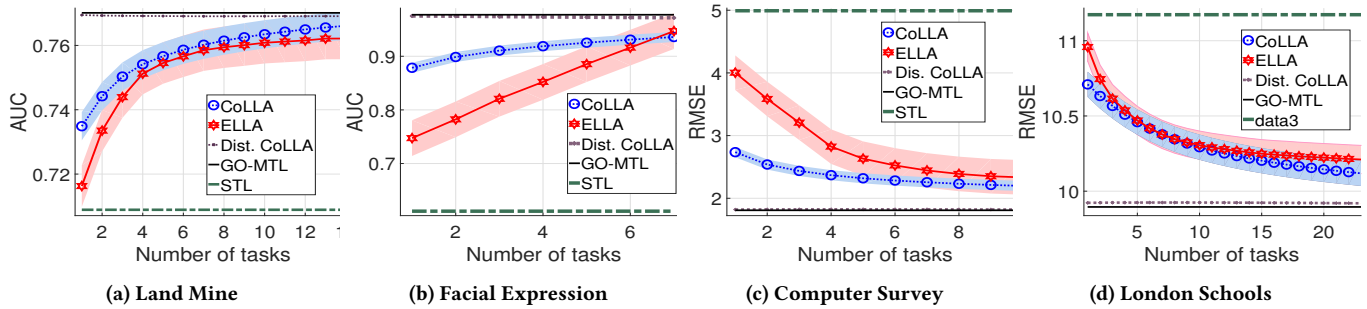


Figure 1: Performance of distributed (dotted lines), centralized (solid), and single-task learning (dashed) algorithms on benchmark datasets. Shaded region shows standard error (Best viewed in color.)

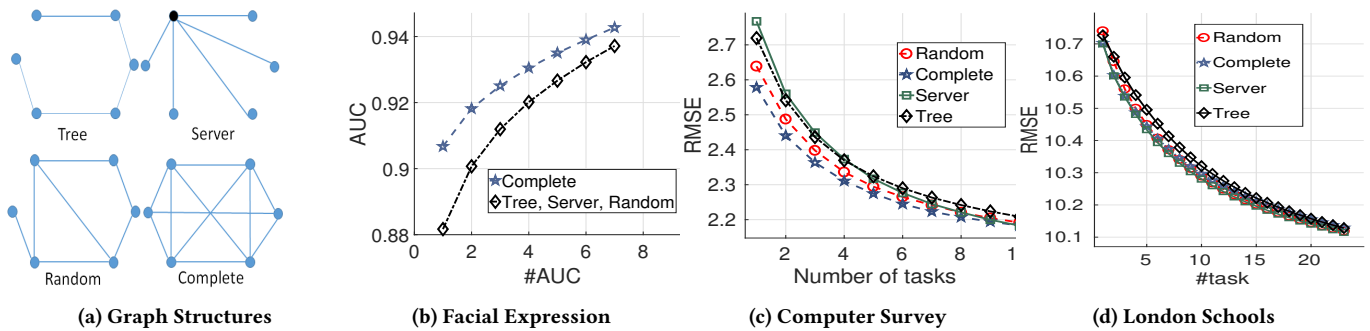


Figure 2: Performance of CoLLA given various graph structures (a) for three data sets (b-d).

of tasks learned by a single agent, but we are reporting the average performance by all agents. We used this visualization, because the benefit of collaborative scheme can be seen on average for all agents.

For the first experiment on CoLLA, we assumed a minimal linearly connected tree which allows for information flow among the agents $\mathcal{E} = \{(i, i + 1) | 1 \leq i \leq N\}$. Figure 1 compares CoLLA against ELLA (which does not use collective learning), GO-MTL, and single-task learning. Note that at each time step t , we report the average performance of online algorithms on all learned tasks up to that instance in the vertical axis. Thus the horizontal axis denotes the number of learned tasks by each agent solely for online learning setting. ELLA can be considered as a special case of CoLLA with an edgeless graph topology (no communication). This would allow us to assess whether consistently positive/ transfer has occurred. A progressive increase in the average performance on the learned tasks demonstrates that positive transfer has occurred and allows plotting learning curves. Moreover, we also performed an offline distributed batch MTL optimization of Eq. (6), i.e. offline CoLLA. For comparison, we plot the learning curves for the online settings and the average asymptotic performance on all tasks for offline settings in the same plot. The shaded regions on the plots denote the standard error for 100 trials.

Figure 1 shows that collaboration among agents expectedly improves lifelong learning, both in terms of learning speed and asymptotic performance, to a level that is not feasible for a single

Method \ Dataset	Dataset			
	LM	LS	CS	FE
CoLLA	6.87	29.62	51.44	40.87
ELLA	6.21	29.30	37.99	38.69
Dist. CoLLA	32.21	37.3	61.71	59.89
GO-MTL	8.63	32.40	61.81	60.17

Table 1: Jumpstart comparison (improvement in percentage) on the Land Mine (LM), London Schools (LS), Computer Survey (CS), and Facial Expression (FE) datasets

lifelong learning agent. Also, performance of offline CoLLA is comparable with GO-MTL, demonstrating that our algorithm can be used as a distributed MTL algorithm. As expected, both CoLLA and ELLA lead to the same asymptotic performance because they solve the same optimization problem. These results demonstrate the effectiveness of our algorithm for both offline and online optimization settings. We also measured the improvement in the initial performance on a new task due to transfer (the “jumpstart” [27]) in Table 1, highlighting CoLLA’s effectiveness in collaboratively learning knowledge bases suitable for transfer. This result accords with intuition because collaboration is most effective in learning initial tasks.

We conducted a second set of experiments to study the effect of the communication mode (i.e., the graph structure) on distributed

lifelong learning. We performed experiments on four graph structures visualized in Figure 2a: tree, server (star graph), complete, and random. The server graph structure connects all agents (slave agents) through a central server (master agent) (depicted in black in the figure), and the random graph was formed by randomly selected half of the edges of a complete graph while still ensuring that the resulting graph was connected. Note that some of these structures coincide when the network is small (for this reason, results on the land mine dataset are not presented for this second experiment). Performance results for these structures on the London schools, computer survey, and facial expression recognition datasets are presented in Figures 2b–2d. Note that for facial recognition data set, results for the only two possible structures are presented. From these figures, we can roughly conclude that for network structures with more edges, the learning rate is faster. Intuitively, this empirical result signals that more communication and collaboration between the agents can increase learning speed.

7 CONCLUSION

We proposed a distributed optimization algorithm for enabling collective multi-agent lifelong learning. Collaboration among the agents not only improves the asymptotic performance on the learned tasks, but enables the agent to learn faster (i.e. using less data to reach a specific performance threshold). Our experiments demonstrated that the proposed algorithm outperforms other alternatives on a variety of MTL regression and classification problems. Extending the proposed framework to a network of asynchronous agents with dynamic links is a potential future direction to improve the applicability of the algorithm on real world problems.

REFERENCES

- [1] A. Argyriou, T. Evgeniou, and M. Pontil. 2008. Convex Multi-Task Feature Learning. *73*, 3 (2008), 243–272.
- [2] I. M. Baytas, M. Yan, A. K. Jain, and J. Zhou. 2016. Asynchronous Multi-task Learning. In *IEEE 16th International Conference on Data Mining (ICDM)*. 11–20.
- [3] S. Ben-David, J. Blitzer, K. Crammer, and F. Pereira. 2007. Analysis of Representations for Domain Adaptation. In *Neural Information Processing Systems*. 137–144.
- [4] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. 2011. Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers. *Foundations and Trends® in Machine Learning* 3, 1 (2011), 1–122.
- [5] R. Caruana. 1997. Multitask Learning. *Machine Learning* 28 (1997), 41–75.
- [6] V. Cevher, S. Becker, and M. Schmidt. 2014. Convex optimization for big data: Scalable, randomized, and parallel algorithms for big data analytics. *IEEE Signal Processing Magazine* 31, 5 (2014), 32–43.
- [7] J. Chen, C. Richard, and A. H. Sayed. 2014. Multitask diffusion adaptation over networks. *IEEE Transactions on Signal Processing* 62, 16 (2014), 4129–4144.
- [8] J. Chen and A. H. Sayed. 2012. Diffusion adaptation strategies for distributed optimization and learning over networks. *IEEE Transactions on Signal Processing* 60, 8 (2012), 4289–4305.
- [9] S. El Bsat, H. Bou-Ammar, and M. E. Taylor. 2017. Scalable Multitask Policy Gradient Reinforcement Learning. In *AAAI*. 1847–1853.
- [10] D. Han and X. Yuan. 2012. A note on the Alternating Direction Method of Multipliers. *Journal of Optimization Theory and Applications* 155, 1 (2012), 227–238.
- [11] J. He and R. Lawrence. 2011. A Graph-based Framework for Multi-task Multi-view Learning. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*. 25–32.
- [12] T. Jansen and R. Wiegand. 2003. Exploring the explorative advantage of the cooperative coevolutionary (1+1) EA. In *Genetic and Evolutionary Computation Conference*. Springer, 310–321.
- [13] X. Jin, P. Luo, F. Zhuang, J. He, and Qing. He. 2015. Collaborating etween Local and Global Learning for Distributed Online Multiple Tasks. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*. ACM, 113–122.
- [14] A. B. Kao, N. Miller, C. Torney, A. Hartnett, and I. D. Couzin. 2014. Collective learning and optimal consensus decisions in social animal groups. *PLoS computational biology* 10, 8 (2014), e1003762.
- [15] A. Kumar and H. Daume III. 2012. Learning task grouping and overlap in multi-task learning. In *Proceedings of the 29th International Conference on Machine Learning*. 1383–1390.
- [16] P. J. Lenk, W. S. DeSarbo, P. E. Green, and M. R. Young. 1996. Hierarchical Bayes conjoint analysis: Recovery of partworth heterogeneity from reduced experimental designs. *Marketing Science* 15, 2 (1996), 173–191.
- [17] M. Li, D. G. Andersen, A. J. Smola, and K. Yu. 2014. Communication efficient distributed machine learning with the parameter server. In *Advances in Neural Information Processing Systems*. 19–27.
- [18] N. Li and J. R. Marden. 2013. Designing games for distributed optimization. *IEEE Journal of Selected Topics in Signal Processing* 7, 2 (2013), 230–242.
- [19] S. Liu, S. J. Pan, and Q. Ho. 2017. Distributed multi-task relationship learning. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 937–946.
- [20] J. Mairal, F. Bach, J. Ponce, and G. Sapiro. 2010. Online learning for matrix factorization and sparse coding. *Journal of Machine Learning Research* 11, Jan, 19–60.
- [21] D. Mateos-Núñez and J. Cortés. 2015. Distributed optimization for multi-task learning via nuclear-norm approximation. The authors are with the Department of Mechanical and Aerospace Engineering, University of California, San Diego, USA. *IFAC-PapersOnLine* 48, 22 (2015), 64–69.
- [22] A. Maurer, M. Pontil, and B. Romera-Paredes. 2013. Sparse coding for multitask and transfer learning. *International Conference on Machine Learning* 28 (2013), 343–351.
- [23] J. F. C. Mota, J. M. F. Xavier, P. M. Q. Aguiar, and M. Puschel. 2013. D-ADMM: A communication-efficient distributed algorithm for separable optimization. *IEEE Transactions on Signal Processing* 61, 10 (2013), 2718–2723.
- [24] S. J. Pan and Q. Yang. 2010. A Survey on Transfer Learning. *IEEE Trans. on Knowledge and Data Engineering* 22, 10 (2010).
- [25] S. Parameswaran and K. Q. Weinberger. 2010. Large margin multi-task metric learning. In *Advances in neural information processing systems*. 1867–1875.
- [26] P. Ruvolo and E. Eaton. 2013. ELLA: An efficient lifelong learning algorithm. *International Conference on Machine Learning* 28 (2013), 507–515.
- [27] M. E. Taylor and P. Stone. 2009. Transfer learning for reinforcement learning domains: A survey. *The Journal of Machine Learning Research* 10 (2009), 1633–1685.
- [28] S. Thrun. 1996. Is learning the n-th thing any easier than learning the first? *Advances in Neural Information Processing Systems* (1996), 640–646. <https://doi.org/10.1.1.44.2898>
- [29] R. J. Tibshirani et al. 2013. The lasso problem and uniqueness. *Electronic Journal of Statistics* 7 (2013), 1456–1490.
- [30] M. F. Valstar, B. Jiang, M. Mehu, M. Pantic, and K. Scherer. 2011. The first facial expression recognition and analysis challenge. In *2011 IEEE International Conference on Automatic Face & Gesture Recognition and Workshops (FG 2011)*. IEEE, 921–926.
- [31] B. Wang and J. Pineau. 2016. Generalized Dictionary for Multitask Learning with Boosting. In *IJCAI*. 2097–2103.
- [32] J. Wang, M. Kolar, and N. Srebro. 2016. Distributed multi-task learning. (2016), 751–760.
- [33] Q. Wang, L. Ruan, and L. Si. 2014. Adaptive Knowledge Transfer for Multiple Instance Learning in Image Classification. In *AAAI*. 1334–1340.
- [34] L. Xie, I. M. Baytas, K. Lin, and J. Zhou. 2017. Privacy-Preserving Distributed Multi-Task Learning with Asynchronous Updates. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 1195–1204.
- [35] E. P. Xing, Q. Ho, W. Dai, J. K. Kim, J. Wei, S. Lee, X. Zheng, P. Xie, A. Kumar, and Y. Yu. 2015. Petuum: A new platform for distributed machine learning on big data. *IEEE Transactions on Big Data* 1, 2 (2015), 49–67.
- [36] Y. Xue, X. Liao, L. Carin, and B. Krishnapuram. 2007. Multi-Task Learning for Classification with Dirichlet Process Priors. *Journal of Machine Learning Research* 8 (2007), 35–63.
- [37] F. Yan, S. Sundaram, S. Vishwanathan, and Y. Qi. 2013. Distributed autonomous online learning: Regrets and intrinsic privacy-preserving properties. *IEEE Transactions on Knowledge and Data Engineering* 25, 11 (2013), 2483–2493.
- [38] D. Zhang, D. Shen, Alzheimer’s Disease Neuroimaging Initiative, et al. 2012. Multi-modal multi-task learning for joint prediction of multiple regression and classification variables in Alzheimer’s disease. *NeuroImage* 59, 2 (2012), 895–907.
- [39] F. Zhang, J. Cao, W. Tan, S. U. Khan, K. Li, and A. Y. Zomaya. 2014. Evolutionary scheduling of dynamic multitasking workloads for big-data analytics in elastic cloud. *IEEE Transactions on Emerging Topics in Computing* 2, 3 (2014), 338–351.
- [40] R. Zhang and J. Kwok. 2014. Asynchronous distributed ADMM for consensus optimization. In *International Conference on Machine Learning*. 1701–1709.