

Full Terms & Conditions of access and use can be found at  
<https://www.tandfonline.com/action/journalInformation?journalCode=uiie21>



# Federated data analytics: A study on linear models

Xubo Yue<sup>a</sup>, Raed Al Kontar<sup>a</sup>, and Ana María Estrada Gómez<sup>b</sup>

<sup>a</sup>Industrial & Operations Engineering, University of Michigan, Ann Arbor, MI, USA; <sup>b</sup>School of Industrial Engineering, Purdue University, West Lafayette, IN, USA

## ABSTRACT

As edge devices become increasingly powerful, data analytics are gradually moving from a centralized to a decentralized regime where edge computing resources are exploited to process more of the data locally. This regime of analytics is coined as Federated Data Analytics (FDA). Despite the recent success stories of FDA, most literature focuses exclusively on deep neural networks. In this work, we take a step back to develop an FDA treatment for one of the most fundamental statistical models: linear regression. Our treatment is built upon hierarchical modeling that allows borrowing strength across multiple groups. To this end, we propose two federated hierarchical model structures that provide a shared representation across devices to facilitate information sharing. Notably, our proposed frameworks are capable of providing uncertainty quantification, variable selection, hypothesis testing, and fast adaptation to new unseen data. We validate our methods on a range of real-life applications, including condition monitoring for aircraft engines. The results show that our FDA treatment for linear models can serve as a competing benchmark model for the future development of federated algorithms.

## ARTICLE HISTORY

Received 13 June 2022

Accepted 6 December 2022

## KEYWORDS

Federated data analytics; linear models; hierarchical model; uncertainty quantification; variable selection; hypothesis testing; fast adaptation; engineering applications

## 1. Introduction

The sheer amount of data collected nowadays is beginning to overwhelm traditional centralized data analytics regimes where data from the edge is continuously uploaded to a central server to be processed. Excessive communication traffic from data upload, significant central server storage needs, energy expenditures from centralized learning of big data models, and privacy concerns from sharing raw data are becoming critical challenges in centralized systems. Statista, a German company specializing in market and consumer data, has predicted that, by 2024, data produced on edge devices (e.g., cell phone data, self-driving vehicle data) will reach more than hundreds of zettabytes while the global central servers only have 10.4 zettabytes of storage (Morell and Alba, 2022). Transmitting such a vast amount of edge data into a central server is infeasible. Adding to that, training a model with moderately large datasets results in significant budget costs and carbon emissions (Patterson *et al.*, 2021). Furthermore, data-sharing comes with serious privacy concerns. According to Lawson *et al.* (2015), Canadian drivers who refused to enroll in the automotive telematics program demanded that their personal driving data (e.g., behavior, location, web-browsing history) should be respected by vehicle companies and that they be given control over the data collection process. These debates over data protection standards have not faded away over the past decade.

Fortunately, the Internet of Things (IoT) is undergoing a new revolution in which the computing power of edge devices is tremendously increasing (Hassan *et al.*, 2018). AI

Chips such as general-purpose chips (GPUs), semi-customized chips (FGPAs), and fully-customized chips (ASICs) are becoming readily available across many applications (Blanco-Filgueira *et al.*, 2019; Rahman and Hossain, 2021; Zhu *et al.*, 2021). Such AI chips are able to process a vast amount of data locally and provide timely responses and decisions (Shi *et al.*, 2016). For instance, the autonomous vehicle company PerceptIn has released a real-time edge computing system, DragonFly+, that is three times more power efficient and delivers three to five times of the computing power of an Nvidia Tx1 and an Intel Core i7 processor (Liu *et al.*, 2019). Another notable example is Tesla's autopilot system that has computing power on the car itself comparable to hundreds of MacBook pros (CleanTechnica, 2021). As a consequence, the traditional IoT is on the verge of shifting to a decentralized framework recently termed the Internet of Federated Things (IoFT) (Kontar *et al.*, 2021) in which some of the data processing is deferred to the edge. In this future, the central server only acts as an orchestrator of the learning process and an integration point of model updates from different devices, rather than the central location where all data is processed. Indeed, IoFT is slowly infiltrating various fields such as manufacturing, transportation, and energy systems (Kontar *et al.*, 2021).

The underlying data analytics framework in IoFT is Federated Data Analytics (FDA), where edge devices exploit their own computation power to collaboratively extract knowledge and build smart analytics while keeping their personal data stored locally. Consequently, edge devices no

longer need to upload their data to the cloud (or server), and, in turn, the cloud does not need to store that immense amount of data. As such, FDA resolves many of the aforementioned drawbacks of the centralized computing system and sets forth many intrinsic advantages, including privacy-preserving and reducing storage/computation/communication costs, among many others.

In spite of some recent advances in FDA, most, if not all, literature focuses on deep neural networks (Li, Sahu, Talwalkar, *et al.*, 2020) (learned via first-order methods). To date, very few papers have delivered federated treatments of traditional statistical models. Perhaps the closest field where statistical models were investigated is Distributed Learning (DL) (Jordan *et al.*, 2019), yet DL and FDA have several fundamental differences. Despite the terminology “distributed”, DL is still a centralized computation approach where different computing nodes operate on all data (Fan *et al.*, 2021). These nodes communicate often, observe each other’s data, and can operate on different data partitions. The underlying philosophy for DL is “divide-and-conquer” where data is divided across the nodes (often dynamically), and then the nodes collaborate to “conquer” (learn) a single model. As a notable example, Zhang *et al.* (2015) propose a DL algorithm that solves a ridge regression problem. The basic idea of this paper is as follows: a server first evenly divides a set of data into  $m$  disjoint sets and assigns each set to a different node. Each node then solves a ridge regression problem and sends the optimal solutions back to the server. The server then aggregates local estimations. As a result, this approach returns a single global model. In contrast, in FDA, data resides at the edge and cannot be shuffled, randomized, or divided. Therefore, edge devices cannot see each others’ data and data partitions for FDA are fixed and often heterogeneous. In addition, devices have datasets with unique features as they correspond to different clients, components or systems (e.g., cars). Therefore, in FDA we cannot divide the data and there is often no single model to conquer, rather, our goal is to borrow strength across edge devices to improve inference and prediction.

In this work, we take a step back and move out of the regime of deep neural networks to study one of the most fundamental statistical models: Linear Regression (LR). Indeed, linear models may facilitate hypothesis testing, uncertainty quantification, variable selection, deriving engineering insight, and establishing a baseline with which to compare other models. Needless to say, in reality, many real applications can be sufficiently characterized by linear models (Liu *et al.*, 2013; Si *et al.*, 2017; Li *et al.*, 2018; Schulz *et al.*, 2020; Arashi *et al.*, 2021; Şahin *et al.*, 2022). In addition, building upon FDA for linear models, one may develop approaches for more complex derivatives such as logistic regression, mixed-effects, and kernel methods.

To this end, we exploit the properties and structure of linear models and develop an FDA treatment for LR with Gaussian noise, entitled FedLin. Our treatment is built upon hierarchical models (HMs), which allow borrowing statistical knowledge across groups (i.e., devices or clients in FDA). Specifically, we propose two federated HM structures

that provide a shared representation across devices to facilitate information transfer. The first structure establishes a shared representation defined through a structural prior over concatenated device parameters. The second structure is based on the assumption that all device parameters are generated from the same underlying distribution. This allows uncertainty quantification and, consequently, a Bayesian treatment for variable selection in FedLin. Our methods are validated on a range of real-life problems, including variable selection and condition monitoring. The results highlight the effective performance of our approaches and their ease of implementation, which may help them serve as benchmark models for many future developments of federated statistical algorithms.

**Organization:** The remainder of this article is organized as follows. In Section 2, we conduct a literature review, and introduce the general setting and motivation. In Sections 3 and 4, we present our two model structures and their applications. We validate our proposed models on various simulated and real-life datasets in Sections 5 and 6. Finally, we draw conclusions in Section 7. The codes used in this article, written in R language, can be found at <https://github.com/UMDataScienceLab/Federated-Linear-Models>.

## 2. Background

### 2.1. Literature overview

The idea of FDA was first brought to the forefront of deep learning by McMahan *et al.* (2017). In this work, they proposed an FDA algorithm termed federated averaging (FedAvg). The idea of FedAvg is simple: a central server distributes initial deep learning model parameters and the network structure to some selected devices, devices perform local Stochastic Gradient Descent (SGD) steps using their data and send their updated parameters back. The server then takes an average of those parameters to update the global model. This process is termed as one communication round and is iterated several times. Although simple, FedAvg is currently still one of the most competitive benchmark models (Kairouz *et al.*, 2021). To date, some work has been proposed to improve the performance of federated deep learning algorithms. For instance, Yuan and Ma (2020) and Liu *et al.* (2020) and Yuan and Ma (2020) provided several provable techniques to accelerate FedAvg and enable faster convergence. Li, Sanjabi, *et al.* (2019); Du *et al.*, (2021); Yu *et al.* (2020); Yue *et al.* (2021); Du *et al.* (2021) developed variants of FedAvg that ensure uniformly good performance across all devices to achieve fairness. Another line of work aims to develop personalized solutions in FDA as excessive heterogeneity can greatly impact the performance of a single global model (Deng *et al.*, 2020; Fallah *et al.*, 2020; Li *et al.*, 2021). Such approaches usually either follow a train-then-personalize philosophy where a trained global model is fine-tuned on local devices or divide the layers of a neural network into shared and individualized ones (Tan *et al.*, 2022), where devices collaborate to learn the common layers using methods such as FedAvg. From a theoretical perspective, Stich (2018); and Li, Sahu, Zaheer, *et al.* (2020) prove the convergence of FedAvg for convex

functions and homogeneous independent and identically distributed (*i.i.d.*) datasets. Those results are then extended to a non-convex setting by Wang and Joshi (2021). On the other hand, Li, Huang, *et al.* (2019) extend the results of Stich (2018) to the non-*i.i.d.* setting. Furthermore, Shi *et al.* (2021) extend the convergence results to a kernel regime. For a comprehensive overview of the current literature, please refer to Kontar *et al.* (2021).

The major trend of FDA exclusively focuses on neural networks and classification tasks. FDA for statistical models is still scarce. Yue and Kontar (2021) extend the Gaussian process to a federated framework and show that their proposed algorithm can achieve state-of-the-art performance on multi-fidelity modeling problems. Yuan *et al.* (2021) develop a federated composite optimization framework that solves the federated Lasso problem. Tong *et al.* (2020) propose a federated iterative hard thresholding algorithm to tackle the non-convex 0-norm penalized regression problem. The two aforementioned papers mainly formulate penalized regression from a frequentist perspective. In Section 4, we will develop a Bayesian formulation built upon HM for federated penalized regression.

## 2.2. General setting

We start by describing our problem setting. Suppose there exists  $K \geq 2$  edge devices. For device  $k \in [K] := \{1, \dots, K\}$ , the dataset is given as  $\mathbf{D}_k = \{\mathbf{X}_k, \mathbf{Y}_k\}$  with  $N_k$  observations, where  $\mathbf{Y}_k = [y_{k1}, \dots, y_{kN_k}]^\top$  is a  $N_k \times 1$  output vector,  $\mathbf{X}_k = [\mathbf{x}_{k1}, \dots, \mathbf{x}_{kN_k}]$  is a  $d \times N_k$  input matrix and  $\mathbf{x}_{ki} = ([\mathbf{x}_{ki}]_1, \dots, [\mathbf{x}_{ki}]_d)^\top \forall i = 1, \dots, N_k$ . Here,  $d$  is the dimensionality of the input space. In this work, we focus on linear models. More specifically, data on device  $k$  is used to learn a linear model parameterized by  $\theta_k \in \mathbb{R}^d$ . The distribution of  $y_{ki}$  is given as

$$y_{ki} | \mathbf{x}_{ki}, \theta_k \sim \mathcal{N}(\mathbf{x}_{ki}^\top \theta_k, \sigma_k^2), \quad \forall i = 1, \dots, N_k, \quad (1)$$

where  $\sigma_k^2$  is a noise parameter. For the sake of compactness, denote by  $\Theta = (\theta_1, \dots, \theta_K)$  a  $d \times K$  matrix concatenating all device parameters.

Furthermore, we assume that a central server is connected to all edge devices and can facilitate the collaborative model learning process. As such, our goal in FDA is to let devices leverage their commonalities to better learn model parameters  $\Theta$ ; all the while distributing the learning efforts and circumventing the need to share raw data.

## 2.3. FDA and HMs

Since our goal is to borrow strength across devices, the first step is to create a shared representation across individual device models in order to facilitate the inductive transfer of knowledge. Here we adopt the natural hierarchy in FDA where a central server is connected to edge devices and can orchestrate the learning process. Specifically, we assume that individual device parameters  $\theta_k$  at the lower hierarchical level are generated from a set of shared parameters at the higher hierarchical level. Through collaboratively learning these shared parameters in a federated manner, devices

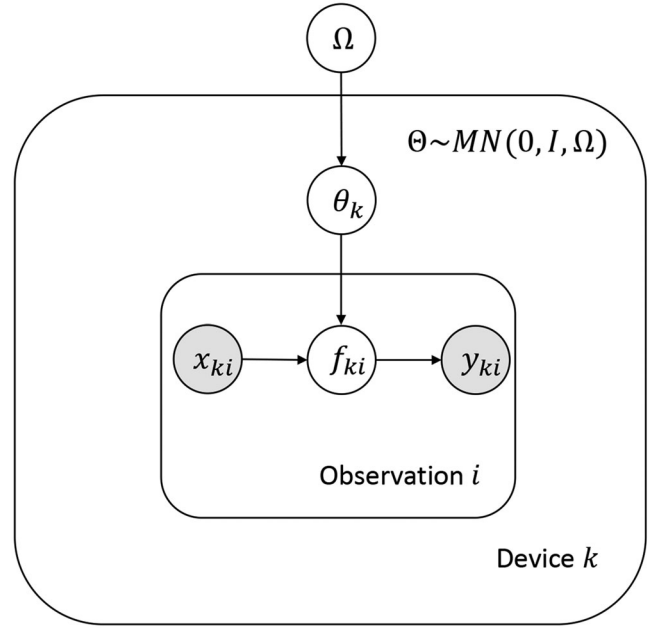


Figure 1. The first HM structure.

induce an update on their personalized parameters  $\theta_k$  that uses information from all other devices.

Two hierarchical structures are proposed. The first defines a joint prior over  $\Theta$  parameterized by a cross-covariance matrix  $\Omega$ . This allows learning a graph that achieves inductive transfer. Whereas, the second HM structure assumes that the  $\theta_k$ 's are sampled from a common distribution (e.g.,  $\theta_k | \phi \sim \mathcal{N}(\mu, \Sigma)$ ). This allows a Bayesian treatment capable of uncertainty quantification as well as learning a global random variable  $\phi$  that can be used to predict on new unseen devices.

We will detail our model formulations, inferences, and applications in the following two sections.

## 3. A shared representation via correlation

In this section, we present our first hierarchical structure (denoted as HM1) that establishes a shared representation by defining a structural prior over  $\Theta$  (Figure 1). This prior is parameterized by  $\Omega$  - a  $K \times K$  cross-covariance matrix. In Figure 1, the matrix  $\Omega$  acts as a graph on the central server that encodes a shared representation among the  $K$  devices and facilitates information sharing. By learning and exploiting the matrix  $\Omega$ , devices can borrow information from each other to improve prediction performance.

Mathematically, we impose a structural prior  $\mathcal{N}(\mathbf{0}, \Omega \otimes \mathbf{I})$  on  $\text{Vec}(\Theta)$ , where  $\text{Vec}(\cdot)$  is a vectorization operation and  $\otimes$  is a Kronecker product. This prior encodes the belief on the underlying distribution that generates components of  $\Theta$ . More specifically,  $\Omega$  is a symmetric matrix whose  $(i, j)$ -th component captures the covariance between device  $i$  and  $j$ . Overall, the aforementioned description translates to the following formulation:

$$\Theta \sim \mathcal{MN}(\mathbf{0}, \mathbf{I}, \Omega), \quad (2)$$

where  $\mathcal{MN}(\mathbf{M}, \mathbf{A}, \mathbf{B})$  denotes a matrix normal distribution with location (mean) parameter  $\mathbf{M}$ , row covariance  $\mathbf{A}$ , and column



covariance  $\mathbf{B}$ . In this prior, the column covariance  $\mathbf{\Omega}$  captures the covariance across devices. Our prior assumes that the covariance across devices is the same for different parameter components. In fact, this is a common practice in the multitask learning literature (Zhang and Yeung, 2012; Ruder, 2017) for multiple reasons: (i) The goal here is to facilitate information transfer between devices and  $\mathbf{\Omega}$  achieves exactly this goal. (ii) This is only a prior and, in reality, it is hard to pre-define the within component covariance. (iii) Posterior computations become rather challenging if the prior (basically a regularization) is complex. As we will show in Sections 5-6, a single  $\mathbf{\Omega}$  is capable of providing excellent performance in several prediction tasks.

By Bayes' rule and incorporating (1) and (2), we can obtain the posterior distribution of  $\mathbf{\Theta}$  as a product of the prior and the likelihood function. By omitting the constant terms and taking the negative logarithm, we can obtain the negative log-likelihood function:

$$\begin{aligned} -\log p(\mathbf{\Theta}|\mathbf{\Omega}, \{\mathbf{Y}_k\}_{k=1}^K) &\propto -\log p(\{\mathbf{Y}_k\}_{k=1}^K|\mathbf{\Theta})p(\mathbf{\Theta}|\mathbf{\Omega}) \\ &= -\log \left( \prod_{k=1}^K p(\mathbf{Y}_k|\mathbf{\Theta}_k) \right) - \log p(\mathbf{\Theta}|\mathbf{\Omega}) \\ &\propto \sum_{k=1}^K \frac{1}{\sigma_k^2} \|\mathbf{Y}_k - \mathbf{X}_k^\top \mathbf{\Theta}_k\|_2^2 + \text{Tr}(\mathbf{\Theta}\mathbf{\Omega}^{-1}\mathbf{\Theta}^\top) \\ &\quad + d \log |\mathbf{\Omega}| := L(\mathbf{\Theta}, \mathbf{\Omega}). \end{aligned}$$

Therefore, our goal is to find the Maximum A Posteriori (MAP) of  $(\mathbf{\Theta}, \mathbf{\Omega})$  that minimizes the negative log-likelihood function, in a federated fashion:

$$(\mathbf{\Theta}^*, \mathbf{\Omega}^*) = \arg \min_{\mathbf{\Theta}, \mathbf{\Omega}} L(\mathbf{\Theta}, \mathbf{\Omega}).$$

To solve this, notice that the derivative with respect to  $\mathbf{\Theta}_k$ , for all  $k$ , is

$$\frac{\partial L(\mathbf{\Theta}, \mathbf{\Omega})}{\partial \mathbf{\Theta}_k} = \frac{-2}{\sigma_k^2} \mathbf{X}_k(\mathbf{Y}_k - \mathbf{X}_k^\top \mathbf{\Theta}_k) + 2 \sum_{i=1}^K \mathbf{\Theta}_i \mathbf{\Omega}_{i,k}^{-1},$$

where  $\mathbf{\Omega}_{i,k}^{-1}$  is defined as the component located in the  $i$ th row and  $k$ th column of  $\mathbf{\Omega}^{-1}$ .

In IoFT, the central server does not have access to datasets  $\mathbf{D} = \{\mathbf{D}_1, \dots, \mathbf{D}_K\}$ , nor do devices have access to each other's datasets. Furthermore, the central server cannot share  $\mathbf{\Omega}$  and  $\mathbf{\Theta}$  with any device, due to the privacy constraint. Therefore, directly running gradient descent using  $\frac{\partial L}{\partial \mathbf{\Theta}_k}$  is not feasible. Yet, by scrutinizing  $\frac{\partial L}{\partial \mathbf{\Theta}_k}$ , one can observe that a gradient update on each  $\mathbf{\Theta}_k$  is split into two gradients. The first term is an update from local data  $\mathbf{D}_k$  whereas the second is a regularization term from all devices based on  $\mathbf{\Omega}$ . Therefore, the local parameter update can be done via the two local steps below

$$\begin{aligned} \text{Stage 1: Multiple local GD or SGD steps } \mathbf{\Theta}_k &\leftarrow \mathbf{\Theta}_k \\ &\quad + 2 \frac{\eta_1}{\sigma_k^2} \mathbf{X}_k(\mathbf{Y}_k - \mathbf{X}_k^\top \mathbf{\Theta}_k). \end{aligned}$$

$$\text{Stage 2: Prior Shrinkage } \mathbf{\Theta}_k \leftarrow \mathbf{\Theta}_k - 2\eta_2 \sum_{i=1}^K \mathbf{\Theta}_i \mathbf{\Omega}_{i,k}^{-1}.$$

At the first stage, device  $k$  runs multiple steps of SGD or Gradient Descent (GD) using the local gradient information

$-\frac{2}{\sigma_k^2} \mathbf{X}_k(\mathbf{Y}_k - \mathbf{X}_k^\top \mathbf{\Theta}_k)$ . To compute this gradient value, one needs to estimate the local variance parameter  $\sigma_k^2$ . Yet, recall that our main goal is to estimate  $\mathbf{\Theta}_k$  by borrowing information from the covariance matrix  $\mathbf{\Omega}$ . Adding to that,  $\mathbf{\Theta}_k$  and  $\sigma_k^2$  are independent. Therefore, it is not necessary to estimate  $\sigma_k^2$  at each local step. Here observe that  $\eta_1$  is a tunable learning rate parameter and we can thus view  $\eta_2 := \frac{\eta_1}{\sigma_k^2}$  as a tuning parameter in stage 1. In other words, we define  $\eta_2$  as the tunable learning rate during the optimization procedure. This circumvents the need to estimate  $\sigma_k^2$  locally. Nevertheless,  $\sigma_k^2$  can be easily estimated from the linear residual term if it is of practitioner's interest. We here note that, since each device  $k$  has a different  $\sigma_k^2$ , it is possible to use device-specific learning rates  $\eta_{2k}$  for all  $k$ . However, it incurs a heavy tuning cost. In this article, we use the same learning rate for each device. As we will show in Section 5-6, this strategy works well.

At the second stage, device  $k$  will then use the aggregated information from all devices  $\sum_{i=1}^K \mathbf{\Theta}_i \mathbf{\Omega}_{i,k}^{-1}$  broadcasted from the central server to update  $\mathbf{\Theta}_k$  by exploiting the covariance matrix  $\mathbf{\Omega}$ . One key notable feature of this updating framework is that the central server only needs to share an aggregated metric  $\sum_{i=1}^K \mathbf{\Theta}_i \mathbf{\Omega}_{i,k}^{-1}$ . This operation is indeed reminiscent of federated averaging and can preserve privacy while allowing devices to borrow strength from each other.

Finally, we will discuss the updating rule of  $\mathbf{\Omega}$  on the central server. The most straightforward way is to take the derivative of  $L(\mathbf{\Theta}, \mathbf{\Omega})$  with respect to  $\mathbf{\Omega}$ . By doing so, we obtain

$$\frac{\partial L(\mathbf{\Theta}, \mathbf{\Omega})}{\partial \mathbf{\Omega}} = -\mathbf{\Omega}^{-1} \mathbf{\Theta}^\top \mathbf{\Theta} \mathbf{\Omega}^{-1} + d \mathbf{\Omega}^{-1} = 0 \Rightarrow \mathbf{\Omega} = \frac{\mathbf{\Theta}^\top \mathbf{\Theta}}{d}.$$

As a result, it is natural to update  $\mathbf{\Omega}$  using this closed-form expression. Here one can view that  $\mathbf{\Theta}^\top \mathbf{\Theta}$  encodes the information of device covariance. Unfortunately, this approach typically faces singularity issues when  $\mathbf{\Theta}^\top \mathbf{\Theta}$  is not a positive definite matrix (e.g., contains zero elements). To resolve this, we propose an updating procedure that prevents an abrupt change in  $\mathbf{\Omega}$  to safeguard against singularity. More specifically, we express the updated  $\mathbf{\Omega}$  as a convex combination between  $\mathbf{\Omega}$  from the previous communication round and the exact updating equation  $\frac{\mathbf{\Theta}^\top \mathbf{\Theta}}{d}$ . That being said, we have

$$\mathbf{\Omega} \leftarrow (1 - \alpha) \mathbf{\Omega} + \frac{\alpha}{d} \mathbf{\Theta}^\top \mathbf{\Theta}.$$

Here,  $\alpha$  is a parameter that controls the change in  $\mathbf{\Omega}$  and  $\frac{\mathbf{\Theta}^\top \mathbf{\Theta}}{d}$  encodes the devices' covariance. A small  $\alpha$  renders a conservative updating rule while a large  $\alpha$  under-weights the importance of the covariance matrix from the previous communication round. When  $\alpha = 1$ , we recover the closed-form updating equation.

Here note that, in the aforementioned framework, the central server selects all devices at each communication round. This scheme is known as full device participation. In reality, however, some local devices are often offline or unwilling to respond due to various reasons. To accommodate this situation, one can sample a subset of devices at each communication round. We term

this scenario as partial device participation. We summarize the detailed algorithm in [Algorithm 1](#).

---

**Algorithm 1:** Improving Device Performance by Exploiting Structural Covariance

---

**Data:** Number of devices  $K$ , number of local iterations  $T$ , set  $\mathcal{S}$  that contains indices of the selected devices, number of communication rounds  $C$ , randomly initialized model parameter  $\Theta$ , initial matrix  $\Omega = I$ , learning rate  $\eta_2$  (selected by grid-search or other tuning methods), proportion  $\alpha = 0.1$  (default), dimension  $d$ .

**for**  $c = 0 : (C - 1)$  **do**

Server broadcasts column of  $\Theta$  (i.e.,  $\theta_k$ ) and  $\sum_{i=1}^K \theta_i \Omega_{i,k}^{-1}$  for all selected  $k$ ;

**for**  $k \in \mathcal{S}$  **do**

**for**  $t = 0 : (T - 1)$  **do**

Device-side: (Sampling Batch) Sample a subset of data  $(\mathbf{X}_k^b, \mathbf{Y}_k^b)$  from  $\mathcal{D}_k$ , where superscript  $b$  means batch;

Device-side: (SGD or GD)  $\theta_k^{(t+1)} = \theta_k^{(t)} + 2\eta_2 \mathbf{X}_k^b (\mathbf{Y}_k^b - \mathbf{X}_k^{b\top} \theta_k^{(t)})$ ;

**end**

Device-side: (Prior Shrinkage)  $\theta_k^{new} \leftarrow \theta_k^{(T)} - 2\eta_2 \sum_{i=1}^K \theta_i \Omega_{i,k}^{-1}$ ;

**end**

Server-side: Combine all  $\theta_k$ , for  $k \notin \mathcal{S}$ , and all  $\theta_k^{new}$ , for  $k \in \mathcal{S}$ , to create a new matrix  $\Theta$ ;

Server-side:  $\Omega \leftarrow (1 - \alpha)\Omega + \frac{\alpha}{d} \Theta^\top \Theta$ .

**end**

Return  $\Theta, \Omega$ .

---

As we will show in our numerical studies, this simple-to-implement algorithm requires very few communication rounds to recover the true parameters and excels at leveraging knowledge across all devices.

#### 4. An HM based on the distribution assumption

So far, we have presented HM1, which exploits the relationship among devices to improve prediction performance. One drawback of [Algorithm 1](#) is that it only returns a point estimate of  $\Theta$ . In practice, it is also desirable to quantify the uncertainty in the parameter estimates. Additionally, the estimated  $\Theta$  and  $\Omega$  cannot provide any borrowable information for new devices, yet the idea of fast adaptation to new unseen data is crucial in many fields such as meta-learning (Vanschoren, 2019). In this section, we will present an alternative model structure that is formulated from a Bayesian perspective to tackle the aforementioned issues.

##### 4.1. Structure and formulation

Our second structure (denoted as HM2) assumes all device parameters are generated from the same underlying distribution. To give a simple example, one can assume  $\theta_k | \phi \sim \mathcal{N}(\mu, \tau I)$  (Figure 2) where  $\phi = (\mu, \tau)$  is a set of a global hyper-parameters

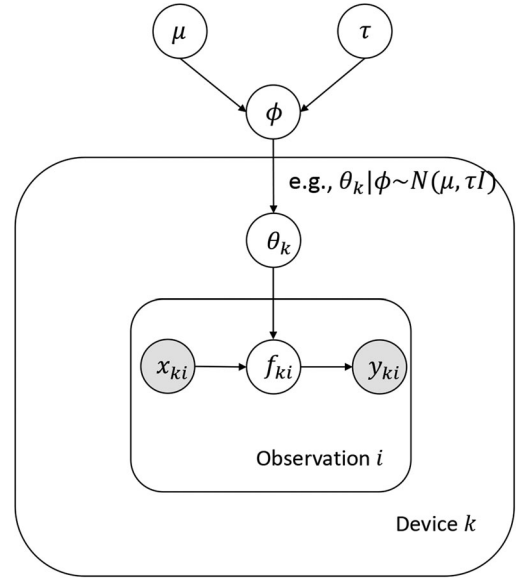


Figure 2. The second HM structure.

on the central server. Here it is critical to note that  $\tau I \in \mathbb{R}^{d \times d}$  is a within covariance matrix and does not denote covariances across  $K$  devices as all  $\theta_k$ 's come from the same underlying distribution. This assignment indicates that  $\{\theta_k\}_{k=1}^K$  are related and generated from the same distribution, yet the degree of model similarity is controlled by the variance parameter  $\tau$ . A small  $|\tau|$  implies all model parameters are similar (i.e., homogeneous) and the hierarchical model is closely related to learning a single common parameter  $\theta$  that fits all devices' data. On the other hand, a large variance  $|\tau|$  incurs more heterogeneity among devices. In the extreme case when  $|\tau| \rightarrow \infty$ , the hierarchical model is equivalent to a separate modeling approach where each device's data is fitted separately (Albert and Hu, 2019).

Now, we formally define the HM2 formulation. From our hierarchical definition and taking a fully Bayesian treatment by placing a prior  $p(\phi)$  on the global hyper-parameters, the joint posterior of  $\{\theta_k\}_{k=1}^K$  and  $\phi$  for HM2 can be written as:

$$\begin{aligned} & p(\theta_1, \dots, \theta_K, \phi | Y_1, \dots, Y_K) \\ & \propto p(Y_1, \dots, Y_K | \theta_1, \dots, \theta_K, \phi) p(\theta_1, \dots, \theta_K, \phi) \\ & = p(\phi) \prod_{k=1}^K p(Y_k | \theta_k, \phi) p(\theta_k | \phi). \end{aligned} \quad (3)$$

To further contextualize HM2, we provide an example of a possible formulation.

**Example 4.1.** Assume each device  $k$  fits a linear regression parameterized by  $\theta_k$ . The local dataset is given as  $(\mathbf{X}_k, \mathbf{Y}_k)$  for all  $k$ . Then one possible hierarchical formulation is:

$$\begin{aligned} Y_k | \theta_k, \phi & \sim \mathcal{N}(\mathbf{X}_k^\top \theta_k, \sigma_k^2 I) \\ \theta_k | \phi & \sim \mathcal{N}(\mu, \tau I) \\ \mu & \sim \mathcal{N}(\mathbf{0}, I) \\ \tau & \sim \log \mathcal{N}(0, 1) \\ \phi & = (\mu, \tau). \end{aligned}$$

Clearly, inferring the joint posterior above is very challenging in a federated setting. Yet, in a hierarchical model, if we

know the posterior over the upper hierarchical level  $p(\phi|\{Y_k\}_{k=1}^K)$  (i.e., over global hyper-parameters  $\phi$ ), we can directly use this posterior as a prior to infer the lower level  $p(\theta_k|\{Y_k\}_{k=1}^K)$  parameters locally.

More specifically, by integrating out device parameters from (3), we can derive that

$$\begin{aligned} & \int p(\theta_1, \dots, \theta_K, \phi | Y_1, \dots, Y_K) d\theta_1 \dots d\theta_K \\ &= p(\phi | Y_1, \dots, Y_K) = p(\phi) \prod_{k=1}^K \int p(Y_k | \theta_k, \phi) p(\theta_k | \phi) d\theta_k = p(\phi) \prod_{k=1}^K f_k(\phi), \end{aligned}$$

where  $f_k(\phi) = \int p(Y_k | \theta_k, \phi) p(\theta_k | \phi) d\theta_k$ . As a consequence, our main goal is to collaboratively learn  $p(\phi | Y_1, \dots, Y_K)$  in a federated setting. One key challenge, however, is that the central server does not have any access to any edge dataset  $D_k$  and therefore directly computing  $p(\phi | Y_1, \dots, Y_K)$  is infeasible. For this reason, we resort to a trick based on the approximate inference methods to learn this posterior distribution.

#### 4.2. Federated Bayesian inference - expectation propagation

In this section, we will present the federated inference framework for HM2. Specifically, our goal is to learn the posterior density  $p(\phi | Y_1, \dots, Y_K)$ . In statistics, the most straightforward and popular approaches to do so are Markov Chain Monte Carlo (MCMC) methods. Yet, as will be clear shortly, we argue that sampling methods are not practical in the federated hierarchical setting, due to their sequential nature. Take Gibbs sampling as an example, device 1 needs to sample  $\theta_1$  from the density  $p(\theta_1 | \theta_2, \dots, \theta_K, \phi, Y_1)$  then pass those samples to the central server. The central server then needs to transmit those sampled  $\theta_1$  to device 2, and device 2 will sample from  $p(\theta_2 | \theta_1, \dots, \theta_K, \phi, Y_2)$ . It can be seen that this sequential nature of MCMC significantly increases the communication cost and also slows down the federated optimization process when the total number of devices is large. Even if one can smartly parallelize the sampling process, the number of MCMC samples obtained locally will be large if the dimension is high, due to the curse of dimensionality (Jordan *et al.*, 2019).

To resolve the aforementioned issues, we resort to Expectation Propagation (EP) (Minka, 2001) to approximate the posterior distribution. EP is one of the most widely-used algorithms for computing an approximate posterior distribution (Minka, 2013; Vehtari *et al.*, 2020). Here, we first briefly introduce the idea of EP in a centralized regime. Consider a posterior distribution with independent data points

$$\pi(\phi) := p(\phi | Y) \propto p(\phi) \prod_{i=1}^N p(y_i | \phi)$$

where  $Y = (y_1, \dots, y_N)^\top$  is the data vector. EP approximates  $\pi(\phi)$  by a density  $q(\phi)$  such that

$$q(\phi) = p(\phi) \prod_{i=1}^N q_i(\phi).$$

Intuitively, EP uses  $q_i(\phi)$  to approximate  $p(y_i | \phi)$ , for all  $i$ . The most common choice is the normal density (Vehtari *et al.*, 2020). To achieve this goal, at each iteration, EP first takes an approximation factor  $q_i(\phi)$  out from the current  $q(\phi)$  and replaces it with the true factor  $p(y_i | \phi)$ . This step yields a new density  $q^{new}(\phi)$ . This resulting new density can be used as an updated approximated posterior. This step is iterated over all  $i$  till convergence. Please, refer to Barthelme (2016) for a comprehensive summary of EP. It can be seen that EP can be naturally extended to FDA where each device can be viewed as an independent “data point”. In the following paragraphs, we will detail the federated extension of EP.

The main idea is to approximate terms  $f_k(\phi)$  by a local device approximation function  $q_k(\phi)$  for all  $k = 1, \dots, K$ . More specifically, we have

$$p(\phi | Y_1, \dots, Y_K) \approx p(\phi) \prod_{k=1}^K q_k(\phi) := q(\phi). \quad (4)$$

Using the framework of EP, we gradually update  $q(\phi)$  by iteratively renewing  $q_k(\phi)$  at each device  $k$ . During each communication round, given the estimated  $q(\phi)$  broadcasted from the server, device  $k$  first computes the cavity distribution

$$q_{-k}(\phi) \propto \frac{q(\phi)}{q_k(\phi)} \quad (5)$$

and the tilted distribution

$$q_{\setminus k}(\phi) \propto f_k(\phi) q_{-k}(\phi). \quad (6)$$

It then computes the updated posterior approximation such that

$$q^{new}(\phi) \approx q_{\setminus k}(\phi). \quad (7)$$

Intuitively, the cavity distribution  $q_{-k}(\phi)$  removes the impact of the old  $q_k(\phi)$  from the approximated posterior density  $q(\phi)$  and the tilted distribution adds the true target density  $f_k(\phi)$  to  $q_{-k}(\phi)$ . As a result, we use the tilted distribution as an updated approximation to the posterior density of  $\phi$ . This step is typically done through a sampling method and is distribution-dependent. We will detail this approximation procedure in Section 4.2.2.

Afterward, device  $k$  calculates the change in its local approximation by

$$\Delta q_k(\phi) = \frac{q^{new}(\phi)}{q(\phi)}. \quad (8)$$

Instead of sending  $q^{new}$  back to the server, we calculate the change in the global posterior imposed by device  $k$  via (8) and sends  $\Delta q_k(\phi)$  to the central server. The server aggregates all device approximations by

$$q(\phi) \leftarrow q(\phi) \prod_{k=1}^K \Delta q_k(\phi). \quad (9)$$

We summarize the EP algorithm in Algorithm 2.

**Algorithm 2:** The Federated Expectation Propagation Algorithm

---

**Data:** number of devices  $K$ , set  $\mathcal{S}$  that contains indices of the selected devices, number of communication rounds  $C$ , initial approximation  $\{q_k(\phi)\}_{k=1}^K$ , prior  $p(\phi)$ .

**for**  $c = 0 : (C - 1)$  **do**

  Server broadcasts  $q(\phi)$ ;

**for**  $k \in \mathcal{S}$  **do**

    Device-side: Calculate the cavity distribution  $q_{-k}(\phi)$  using Eq. (5);

    Device-side: Calculate the tilted distribution  $q_{\setminus k}(\phi)$  using Eq. (6);

    Device-side: Get new  $q(\phi)$  from the tilted distribution using Eq. (7);

    Device-side: Calculate  $\Delta q_k(\phi)$  using Eq. (8) and update local  $q_k(\phi)$ ;

    Device-side: Send  $\Delta q_k(\phi)$  to the central server;

**end**

  Server-side: Update  $q(\phi)$  using Eq. (9);

**end**

Return  $q(\phi)$ .

---

**4.2.1. Posterior of device parameters**

Once we obtain  $q(\phi)$  that approximates  $p(\phi|Y_1, \dots, Y_K)$ , we can further estimate the posterior of device parameters. Specifically, given a device  $k$ ,

$$\begin{aligned}
 p(\theta_k | \{Y_k\}_{k=1}^K) &= \int_{\phi} \int_{\theta_j, j \neq k} p(\theta_k, \phi | \{Y_k\}_{k=1}^K) d\theta_j d\phi \\
 &\propto \int p(\phi) p(Y_k | \theta_k, \phi) p(\theta_k | \phi) \\
 &\quad \prod_{j \neq k} \int p(Y_j | \theta_j, \phi) p(\theta_j | \phi) d\theta_j d\phi \\
 &\approx \int q_{-k}(\phi) p(Y_k | \theta_k, \phi) p(\theta_k | \phi) d\phi.
 \end{aligned}$$

As a consequence, we can use the posterior of  $\theta_k$  to quantify uncertainties or conduct hypothesis testing. The posterior samples from  $p(\theta_k | \{Y_k\}_{k=1}^K)$  can be obtained by off-the-shelf posterior sampling techniques (see `mcmc` package in R or `NumPyro` library in Python). Here we provide a simpler sampling trick. In the above equation, if we ignore the integral, we can obtain the joint posterior

$$p(\theta_k, \phi | \{Y_k\}_{k=1}^K) \approx q_{-k}(\phi) p(Y_k | \theta_k, \phi) p(\theta_k | \phi). \quad (10)$$

As a result, one can ignore the integration and use sampling methods to jointly sample  $(\theta_k, \phi)$  from (10) and discard  $\phi$ . Now, given  $M$  samples  $\{\theta_{ki}\}_{i=1}^M$ , we can readily use the samples to estimate moments, coverage probability, and do hypothesis tests.

Here we should note that the estimated posterior density  $q(\phi)$  encodes crucial information across all devices (recall the explanation in Section 4.1). One can exploit this information for a new device to achieve fast adaptation. For example, we can treat the posterior mean of  $\phi$  as an initial

model parameter for a new device. This idea is similar to meta-learning (Vanschoren, 2019), where one tries to learn a global model that can quickly adapt to a new task.

**4.2.2. Normal approximation**

In practice, it is common to model  $p(\phi)$  and  $q_k(\phi)$ ,  $\forall k$  as normal densities. This is due to a very useful property of normal random variables.

**Lemma 1.** (Williams and Rasmussen, 2006) Suppose there are two normal random variables (with the same dimension) such that  $\theta_1 \sim \mathcal{N}(\mu_1, \Sigma_1)$  and  $\theta_2 \sim \mathcal{N}(\mu_2, \Sigma_2)$ . Let  $r_i = \Sigma_i^{-1} \mu_i$ ,  $Q_i = \Sigma_i^{-1}$  for  $i = 1, 2$ . Define  $p(\theta_+) = p(\theta_1)p(\theta_2)$  and  $p(\theta_-) = \frac{p(\theta_1)}{p(\theta_2)}$ . We have that

$$\begin{aligned}
 \theta_+ &\sim \mathcal{N}(r_1 + r_2, Q_1 + Q_2) \\
 \theta_- &\sim \mathcal{N}(r_1 - r_2, Q_1 - Q_2).
 \end{aligned}$$

Lemma 1 states that the product of two Gaussian densities gives another unnormalized Gaussian density. We use  $\theta_+$  to represent this new Gaussian random variable. Similarly, the quotient of two Gaussian densities gives an unnormalized Gaussian density and we use  $\theta_-$  to represent this new random variable.

Using Lemma 1, one can efficiently implement the EP algorithm, as all components in Algorithm 2 can be computed in closed forms. Here, we detail the implementation technique. We model the prior of  $\phi$  as a multivariate normal random variable with mean  $\mu_0$  and variance  $\Sigma_0$ . We also assume  $q_k(\phi)$  has a normal density parameterized by  $\mu_k, \Sigma_k$ , for all  $k$ . If the support of some components in  $\phi$  does not lie in  $\mathbb{R}$ , one can always perform a logarithmic or logistic transformation to those components. By Gaussian properties, (4) can be computed in closed-form such that  $q(\phi)$  has a normal density parametrized by mean  $r_0 + \sum_{k=1}^K r_k$  and variance  $Q_0 + \sum_{k=1}^K Q_k$ , where  $r_j = \Sigma_j^{-1} \mu_j$  and  $Q_j = \Sigma_j^{-1}$ , for all  $j = 0, 1, \dots, K$ . Similarly, by Lemma 1, the cavity distribution in (5) can be computed in a closed-form by a subtraction operation. This results in  $r_{-k} := r - r_k$ ,  $Q_{-k} := Q - Q_k$ . Compared to sampling approaches, one key advantage of EP is that the computation and communication steps are simple and efficient. The central server and devices only need to transmit the mean vector and variance matrix to perform model updating and aggregation. Notably, during each communication round, there is no need to estimate or transmit normalizing constants. We detail this idea in Algorithm 3. In Algorithm 3, one needs to compute the tilted distribution and obtain  $r_k^{new}, Q_k^{new}$ . These steps correspond to (6) and (7). However, as  $f_k(\phi)$  may not be a normal density, the resulting tilted distribution is not normal. Therefore, we need to use a normal distribution to approximate the tilted distribution. To do so, we can run a set of simulation draws (using any sampling method coded in R or Python libraries) from  $f_k(\phi)q_{-k}(\phi)$  and estimate the mean and covariance of those draws. We set the resulting mean to be  $\mu_{\setminus k}$  and the resulting covariance to be  $\Sigma_{\setminus k}$ .



---

**Algorithm 3:** The Federated Expectation Propagation Algorithm using Normal Approximation

---

**Data:** number of devices  $K$ , set  $\mathcal{S}$  that contains indices of the selected devices, number of communication rounds  $C$ , initial approximation  $\{\mathbf{r}_k, \mathbf{Q}_k\}_{k=1}^K$ , prior  $\mathbf{r}_0, \mathbf{Q}_0$ , initial posterior parameters  $\mathbf{r} = \mathbf{r}_0 + \sum_{k=1}^K \mathbf{r}_k, \mathbf{Q} = \mathbf{Q}_0 + \sum_{k=1}^K \mathbf{Q}_k$ .

**for**  $c = 0 : (C - 1)$  **do**  
 Server broadcasts  $\mathbf{r}, \mathbf{Q}$ ;  
**for**  $k \in \mathcal{S}$  **do**  
   Device-side: Calculate the cavity distribution with parameters  $\mathbf{r}_{-k} := \mathbf{r} - \mathbf{r}_k, \mathbf{Q}_{-k} := \mathbf{Q} - \mathbf{Q}_k$ ;  
   Device-side: Calculate the tilted distribution  $\mathbf{r}_{\setminus k}, \mathbf{Q}_{\setminus k}$ ;  
   Device-side: Obtain new  $\mathbf{r}_k^{new}, \mathbf{Q}_k^{new}$ ;  
   Device-side: Calculate  $\Delta \mathbf{r}_k = \mathbf{r}_k^{new} - \mathbf{r}_k, \Delta \mathbf{Q}_k = \mathbf{Q}_k^{new} - \mathbf{Q}_k$ ;  
   Device-side: Send  $\Delta \mathbf{r}_k, \Delta \mathbf{Q}_k$  to the central server;  
**end**  
 Server-side: Update  $\mathbf{r} = \mathbf{r} + \sum_{k \in \mathcal{S}} \Delta \mathbf{r}_k, \mathbf{Q} = \mathbf{Q} + \sum_{k \in \mathcal{S}} \Delta \mathbf{Q}_k$ ;  
**end**  
 Return  $\mathbf{r}, \mathbf{Q}$ .

---

#### 4.3. Federated and penalized regression for variable selection

To move a step further, we ask “is it possible to let devices exploit the shared representation structure to perform variable selection?” Indeed, there are some attempts to tackle this question from a frequentist perspective. Yuan *et al.* (2021) develop a federated composite optimization framework that solves the federated Lasso problem. Tong *et al.* (2020) propose a federated iterative hard thresholding algorithm to tackle non-convex penalized regression. Despite these few efforts in exploring variable selection from a frequentist perspective, no literature exists in the Bayesian setting.

Tibshirani (1996) has shown that a Lasso estimate can be achieved when the regression parameters have *i.i.d.* Laplace priors. Since then, researchers have started to build Bayesian priors for many other penalized regressions such as the elastic net and fused Lasso. Please refer to the work of Van Erp *et al.* (2019) for a detailed literature review. Inspired by the Bayesian interpretation of penalized regressions, we develop a hierarchical structure, based upon HM2, to perform federated variable selection. To proceed, we impose priors on  $\theta_k$ , for all  $k$ , such that

$$\theta_{ki} | \phi \sim \pi(\lambda, \sigma^2), \forall i = 1, \dots, d$$

where  $\phi = (\lambda, \sigma^2)$ ,  $\lambda$  is a regularization parameter and  $\sigma^2$  is a variance parameter. Here,  $\pi(\lambda, \sigma^2)$  is a distribution parameterized by  $\lambda, \sigma^2$ . For instance, if we set  $\pi(\lambda, \sigma^2)$  to be a Laplace distribution with zero mean and  $\frac{\sigma^2}{\lambda}$  diversity, then we recover the Bayesian counterpart of Lasso regression (Tibshirani, 1996). Another example is if we set  $\pi(\lambda, \sigma^2)$  to

be  $\mathcal{N}(0, \frac{\sigma^2}{\lambda})$ , then we recover Ridge regression. There are many possible choices of prior beliefs on  $\sigma^2$  and  $\lambda$ . In this work, we impose log-normal priors on  $\sigma^2$  and  $\lambda$  (Van Erp *et al.*, 2019) and we set  $\phi = (\log \lambda, \log \sigma^2)$ . Our framework can flexibly incorporate other priors such as a non-informative prior on  $\sigma^2$  or a half-Cauchy prior on  $\lambda$ . In this work, we will use a log-normal prior as an illustrative example.

The posterior distribution of  $\theta_k$ , for all  $k$ , and  $\phi$  can be learned by Algorithm 2. Here, one caveat is that, unlike frequentist penalized regressions, the Bayesian methods do not shrink regression coefficients to be exactly zero. As a result, we will calculate the Credible Interval (CI) for each parameter. If the CI of a parameter, say  $\theta_{ki}$ , covers zero, we will exclude this predictor.

## 5. Simulation studies

### 5.1. HM1: Proof of concept using Algorithm 1

**Case I:** We assume that  $K=2$  and generate  $(\theta_1, \theta_2)$  from a matrix-variate normal distribution with zero mean and

$$\mathbf{I} \otimes \mathbf{\Omega} = \mathbf{I} \otimes \begin{bmatrix} 1 & 0.7 \\ 0.7 & 1 \end{bmatrix}$$

covariance. The input space dimension is  $d=5$ . Data on each device are generated from linear models using the generated parameters. We set noise to be 0.05. To demonstrate the benefits of our correlation-based construction in HM1, we create imbalanced sample sizes on the devices. Specifically, device 1 only has  $N_1 = 20$  data points whereas device 2 has  $N_2 = 200$  data points. We train Algorithm 1 with  $C = 30, \eta_2 = 0.01, \alpha = 0.1$  and we set the number of local steps  $T$  to be 20. We compare the performance of Algorithm 1 with a separate modeling approach where each device fits its own model without communication. Specifically, each device runs 600 local SGD steps with learning rate 0.01.

**Case II:** We set  $K=100$  and generate a  $100 \times 100$  positive definite matrix  $\mathbf{\Omega}$  using the R package `clusterGeneration`. We then generate true device parameters based on the matrix  $\mathbf{\Omega}$ . We set  $d=8$  and generate data using the linear models with noise  $\sigma^2 = 0.1$ . For the first 30 devices, we assign 40 data points and for the remaining 70 devices, we assign 275 data points. Overall, we create an imbalanced data generation scenario. Case II can be viewed as a generalization of Case I with more devices. Again, we train models using the same hyperparameters as those in Case I.

**Case III:** We use the same setting as the one in Case II, but all devices have 20 data points (i.e., balanced data generation).

**Case IV:** We increase the sample size on each device to 200 and use the same setting as the one in Case III.

**Performance Evaluation:** Denote by  $(\mathbf{X}_k^*, \mathbf{Y}_k^*)$  the testing dataset on device  $k$  where  $\mathbf{X}_k^* = [x_{k1}^*, \dots, x_{kN_k^*}^*]^\top$  and  $\mathbf{Y}_k^* = [y_{k1}^*, \dots, y_{kN_k^*}^*]$ . The averaged Root-mean-square error (A-

RMSE) across all devices is defined as

$$A - RMSE = \frac{1}{K} \sum_{k=1}^K \sqrt{\frac{\sum_{i=1}^{N_k^*} (f(x_{ki}^*) - y_{ki}^*)^2}{N_k^*}},$$

where  $f_k(x^*) = x^{*\top} \theta_k$ . On each device, we generate 1000 data points using the true device parameters for testing. In Case I, the RMSE is calculated on device 1. In Case II, the A-RMSE is calculated using devices  $k \in \{1, \dots, 30\}$ . Here our goal is to assess the prediction accuracy on devices with scant data. In Cases III/IV, we calculate A-RMSE using all 100 devices. We report our results in Table 1. It can be seen that Algorithm 1 yields much smaller A-RMSE under the imbalanced data scenario. This conveys the importance of borrowing strength from other devices under the FDA framework. In Case III, the local sample size is insufficient to allow each model alone to perform well. However, our FDA can still benefit devices' training by borrowing information from other devices. In case IV, Algorithm 1 does not offer a major improvement, as all local devices have enough data. In this case, doing local training without collaboration should be sufficient.

**Accuracy of Parameter Estimation:** The negative log-likelihood function  $L(\Theta, \Omega)$  is a non-convex function of  $(\Theta, \Omega)$ . To test the impact of the initialization, we conduct a sensitivity analysis below. Denote by  $\hat{\Theta}$  the concatenated estimated device parameters and  $\Theta^*$  the concatenated true data-generating parameters. In Figure 3, we plot  $\frac{\|\hat{\Theta} - \Theta^*\|}{\sqrt{K}}$  versus communication round for Case II and III over 30

independent runs. Each independent run used a different initialized  $\Theta$ . More specifically, we first generate  $d * K$  random numbers from a standard normal distribution. We then create a  $d * K$  matrix  $\Theta$  using these random numbers. It can be seen that Algorithm 1 accurately recovers the true underlying model parameters. Furthermore, it can be observed that Algorithm 1 typically converges within 30–40 communication rounds. We observed that, in all simulations, Algorithm 1 could be trained within 5 seconds on a standard laptop. In conclusion, our proposed algorithm is easy to implement and optimize.

## 5.2. HM2 : Proof of concept

In this section, we test the variable selection performance of HM2. Following the examples in Van Erp *et al.* (2019), we create several simulation cases below:

**Case I:** We set  $K = 10$ ,  $\theta_{true} = (3, 1.5, 0, 0, 2, 0, 0, 0)^\top$  and generate all columns of  $\{X_k\}_{k=1}^K$  from a standard multivariate normal distribution. We then generate  $\{Y_k\}_{k=1}^K$  using  $\theta_{true}$  and  $\{X_k\}_{k=1}^K$  for all  $k$ . We set the noise to 0.05. Each device has 100 data points for training and 1000 data points for testing.

**Case II:** We use the same setting as the one in Case I. The difference is that the first two devices only have 20 data points each, whereas the other devices have 200 data points each. The number of testing data points is 1000.

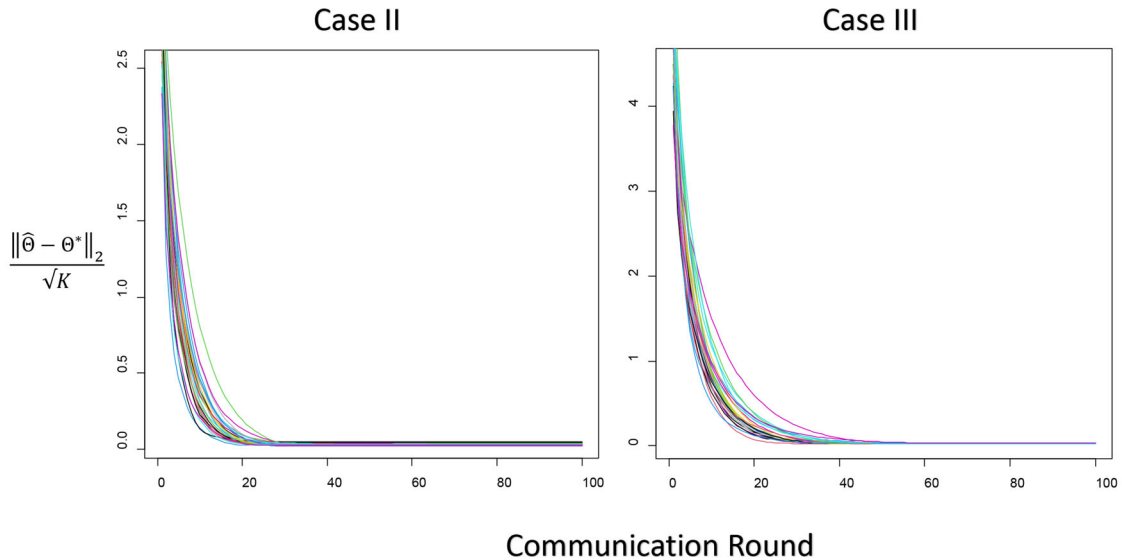
**Case III:** We set  $K = 20$ ,  $\theta_{true} = (\underbrace{3, \dots, 3}_{10}, \underbrace{0, \dots, 0}_{10}, \underbrace{3, \dots, 3}_{10})^\top$ .

Each device has 40 observations for training and 400 observations for testing.

We evaluate the performance of our model based on prediction and variable selection accuracy. The prediction accuracy is evaluated by A-RMSE. Variable selection accuracy is based on the averaged correct and false inclusion rates. To decide whether to include a variable or not, we first calculate a

**Table 1.** The A-RMSE of our proposed model and the separate model over 30 independent runs. We report standard deviations of A-RMSEs in brackets.

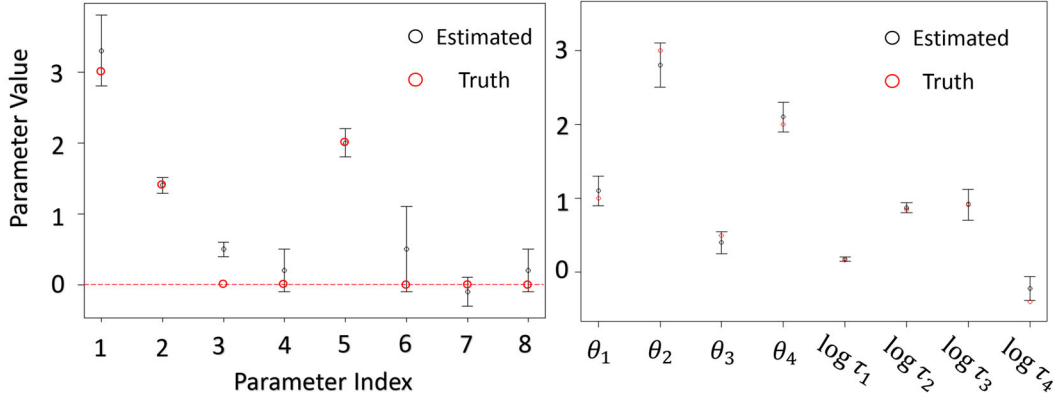
Case	Algorithm 1	Separate
I	0.081 ( $\pm 0.001$ )	0.094 ( $\pm 0.001$ )
II	0.050 ( $\pm 0.000$ )	0.056 ( $\pm 0.001$ )
III	0.044 ( $\pm 0.002$ )	0.072 ( $\pm 0.004$ )
IV	0.035 ( $\pm 0.000$ )	0.035 ( $\pm 0.000$ )



**Figure 3.** Plot of  $\frac{\|\hat{\Theta} - \Theta^*\|}{\sqrt{K}}$  versus communication round. Each color represents one independent run.

**Table 2.** The A-RMSE and averaged correct/false inclusion rates for different federated variable selection methods over 30 experimental runs.

Methods	A-RMSE	Averaged Correct Inclusion Rate	Averaged False Inclusion Rate
Lasso (HM2, Case I)	0.055 ( $\pm 0.001$ )	0.880 ( $\pm 0.003$ )	0.095 ( $\pm 0.001$ )
Lasso (HM2, Case II)	0.062 ( $\pm 0.002$ )	0.875 ( $\pm 0.004$ )	0.101 ( $\pm 0.001$ )
Lasso (HM2, Case III)	0.088 ( $\pm 0.001$ )	0.891 ( $\pm 0.005$ )	0.115 ( $\pm 0.001$ )

**Figure 4.** Plot of parameter estimation and CI.

90% CI for each parameter. If the CI covers zero, we will exclude this predictor. Results are reported in Table 2. It can be seen that our proposed federated variable selection methods can correctly identify more than 85% of effective predictors while maintaining low false inclusion rates.

As mentioned in Section 4, one advantage of HM2 is that it can provide Uncertainty Quantification (UQ) for parameter estimation. We will provide two examples to demonstrate the UQ capability of HM2:

1. We collect the estimated posterior for parameters  $\theta_1$  from an independent run in case I and calculate the mean and 90% credible interval. The resulting plot is presented in Figure 4 (Left). It can be seen that the true parameters are included in the CI generated by HM2.
2. We create  $K=100$  devices and generate device parameter

$$\theta_k | \phi \sim \mathcal{N}(\mu_{true}, \Sigma_{true})$$

$$:= \mathcal{N} \left( \begin{bmatrix} 1 \\ 3 \\ 0.5 \\ 2 \end{bmatrix}, \begin{bmatrix} 1.17 & 0 & 0 & 0 \\ 0 & 2.35 & 0 & 0 \\ 0 & 0 & 2.52 & 0 \\ 0 & 0 & 0 & 0.67 \end{bmatrix} \right)$$

for  $k \in \{1, \dots, 100\}$ . We then use  $\theta_k$  to generate 100 data points for each device  $k$ . Our HM structure can be summarized as follows:

$$\begin{aligned} Y_k | \theta_k, \phi &\sim \mathcal{N}(X_k^\top \theta_k, \sigma_k^2 \mathbf{I}) \\ \theta_k | \phi &\sim \mathcal{N}(\mu, \text{diag}(\tau_1, \dots, \tau_4)) \\ \mu &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \tau_i &\sim \log \mathcal{N}(0, 1), \forall i \\ \phi &= (\mu, \log \tau_1, \dots, \log \tau_4). \end{aligned}$$

We calculate the posterior distribution of  $\phi$  using Algorithm 3 and plot the mean and 90% CI for each

component in Figure 4 (Right). It can be seen that the mean of the posterior of  $\phi$  is close to the truth and the 90% CI also covers the true parameter. This estimated  $q(\phi)$  can be used as an initialization for new devices to achieve fast adaption.

## 6. Real-world case studies

### 6.1. Student performance dataset

This is a public dataset that can be found at <https://archive.ics.uci.edu/ml/index.php>. The dataset contains information on student performance (measured by exam scores) in secondary education of two Portuguese schools, namely, Gabriel Pereira and Mousinho da Silveira. It includes 29 predictors covering gender, grades, demographic, and many other social/school-related features. Detailed information can be found in Cortez and Silva (2008). We treat each school as a “device” (i.e.,  $K=2$ ). On each device, we randomly pick 60% of the students as the training dataset and another 40% of the students as the testing dataset. We create dummy variables for all nominal variables, such as job and guardian. All other numeric variables are standardized to a zero mean and one standard deviation, following the guide in Cortez and Silva (2008). This data processing yields 38 predictors.

Our first goal is to select relevant predictors using our federated penalized regression technique (See Section 4.3). We then use the selected predictors to predict the final exam grades of students. We consider the two most widely-used variable selection methods: Lasso and Ridge. Results are reported in Table 3. The model performance is evaluated based on the RMSE and the number of included predictors.

It can be seen that the variable selection performance of HM2 is consistent with the centralized variable selection method such as Lasso and Ridge regressions. This implies that our framework can serve as a new paradigm for decentralized variable selection problems. Additionally, HM2 also

**Table 3.** The RMSE and number of included predictors for different federated variable selection methods.

Methods	A-RMSE	# included predictors (School 1)	# included predictors (School 2)
Lasso (HM2)	0.825	21	23
Ridge (HM2)	0.817	21	22
Methods	RMSE	# included predictors	
Lasso (Centralized)	0.820	21	
Ridge (Centralized)	0.815	21	

**Table 4.** The A-RMSE of all models over 30 independent runs. We report the standard deviation in the brackets.

Sensor	HM1 ( $\alpha = 0.9$ )	Separate	FedAvg	Ditto	Dis-Ridge
Sensor 2	<b>0.270</b> ( $\pm 0.001$ )	0.299( $\pm 0.003$ )	0.450( $\pm 0.013$ )	0.281( $\pm 0.001$ )	0.552( $\pm 0.002$ )
Sensor 3	<b>0.218</b> ( $\pm 0.002$ )	0.223( $\pm 0.001$ )	0.303( $\pm 0.009$ )	0.220( $\pm 0.001$ )	0.288( $\pm 0.002$ )
Sensor 7	<b>0.369</b> ( $\pm 0.004$ )	0.405( $\pm 0.003$ )	0.628( $\pm 0.011$ )	0.388( $\pm 0.005$ )	0.605( $\pm 0.001$ )
Sensor 8	<b>0.267</b> ( $\pm 0.001$ )	0.307( $\pm 0.001$ )	0.395( $\pm 0.008$ )	0.289( $\pm 0.001$ )	0.390( $\pm 0.003$ )

yields comparable A-RMSEs compared to centralized methods. This demonstrates the advantage of borrowing strength from other devices. However, please note that, in terms of prediction accuracy, federated variable selection can rarely beat the centralized approach as the latter uses more data.

## 6.2. NASA aircraft gas turbine engines

In this case study, we consider condition monitoring data generated from aircraft gas turbine engines using the NASA commercial Modular Aero-Propulsion System Simulation (C-MAPSS) tools. The dataset is available at <https://ti.arc.nasa.gov/tech/dash/groups/pcoe/>. This dataset contains 100 engines. In each engine, 24 sensors are installed to collect time-series degradation signals. For each engine, we treat the first 60% of the time-series observations as the training dataset and the remaining 40% of the signals as the testing dataset. Within the training dataset, we sample 20% of the data as a validation dataset. Our goal is therefore to predict the sensor signal trajectory on each gas turbine engine by training [Algorithm 1](#) using the training dataset. In this scenario, each engine can be viewed as a device (i.e.,  $K=100$ ).

It can be observed that all signal trajectories exhibit polynomial patterns and therefore, many existing works resort to polynomial regression to analyze this dataset (Liu *et al.*, 2013; Song and Liu, 2018). Here we detail the modeling procedure. Given a specific sensor, for all  $k \in \{1, \dots, K\}$ , device  $k$  fits a  $d=6$ th-order polynomial regression in the form of

$$Y_k = \mathbf{X}_k^\top \boldsymbol{\theta}_k + \text{noise},$$

where the  $(d+1) \times N_k$  design matrix  $\mathbf{X}_k$  is in the form of

$$\mathbf{X}_k^\top = \begin{bmatrix} 1 & [x_{k1}]_1 & [x_{k1}]_2^2 & \dots & [x_{k1}]_d^d \\ 1 & [x_{k2}]_1 & [x_{k2}]_2^2 & \dots & [x_{k2}]_d^d \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & [x_{kN_k}]_1 & [x_{kN_k}]_2^2 & \dots & [x_{kN_k}]_d^d \end{bmatrix}.$$

In the above expression,  $Y_k$  represents the signal trajectory for device  $k$  and  $x_{k1}$  represents time. In HM1, device parameters  $\{\boldsymbol{\theta}_k\}_k$  are estimated using [Algorithm 1](#).

We benchmark our proposed model with the following algorithms:

- **FedAvg:** FedAvg is one of the most fundamental and competing benchmark models in the FDA. During each communication round, device  $k$ , for all selected  $k$ , first runs several steps of local SGD and then sends updated parameters  $\boldsymbol{\theta}_k$  back to the server. The server then aggregates those parameters by calculating  $\bar{\boldsymbol{\theta}} = \frac{1}{|\mathcal{S}|} \sum_{k \in \mathcal{S}} \boldsymbol{\theta}_k$ . This step is repeated several times and ultimately, each device will use the global parameter  $\bar{\boldsymbol{\theta}}$  to perform prediction.
- **Dis-Ridge** (Zhang *et al.*, 2015): Dis-Ridge is a distributed ridge regression method. The server first evenly divides a set of data into  $m$  disjoint sets and assigns each set to a different node. Each node then solves a ridge regression problem, by optimizing  $\frac{1}{N_k} \sum_{i=1}^{N_k} (y_{ki} - \mathbf{x}_{ki}^\top \boldsymbol{\theta}_k)^2 + \lambda \|\boldsymbol{\theta}_k\|^2$ , and sends the optimal solution back to the server. The server then aggregates local estimations.
- **Ditto:** Ditto is a personalized FL algorithm. The first stage of Ditto is the same as FedAvg and generates a global parameter  $\bar{\boldsymbol{\theta}}$ . Afterwards, each device  $k$  derives the personalized solution  $\mathbf{v}_k$  by solving a constrained optimization problem  $F_k(\mathbf{v}_k) + \frac{\lambda_{Ditto}}{2} \|\mathbf{v}_k - \bar{\boldsymbol{\theta}}\|_2^2$  where  $F_k(\cdot)$  is the local loss function and  $\lambda_{Ditto}$  is a tuning parameter. The intuition is that each device can run several updating procedures to collect personalized solutions while this solution stays in the vicinity of the shared global model to retain useful information from a global model.
- **Separate:** In Separate, each device simply fits its own linear model without communication.

For all models, we set  $T=20$ ,  $C=100$ , and use grid-search to tune the learning rate (and other model hyperparameters). In [Algorithm 1](#), we set  $\alpha=0.9$ . We report the A-RMSE across all 100 devices in [Table 4](#).

From [Table 4](#), it can be seen that FedAvg and Dis-Ridge consistently yield the worst prediction accuracy as one shared global parameter  $\bar{\boldsymbol{\theta}}$  does not suit all devices, especially in a heterogeneous setting. Personalized approaches such as Ditto circumvent this disadvantage of global models and generate personalized solutions for each



device. Those personalized methods, however, ignore related information amongst devices. Our method, on the other hand, improves the prediction accuracy by exploiting a joint structure for inductive transfer.

## 7. Conclusion

This article proposes a federated treatment for linear regression by adopting a HM approach. We test our proposed framework on a range of simulated and real-world datasets. Despite the simplicity of our linear model framework, it can outperform many state-of-the-art federated algorithms and we argue that it can serve as a competing benchmark model for the future development of federated algorithms. One possible future direction is to extend our framework to generalized linear models such as linear mixed-effect models or to more complicated models such as Gaussian processes and tensor regression. We hope our work will help inspire continued exploration into the world of federated data analytics and its engineering applications.

## Funding

The authors acknowledge the generous support from the NSF CAREER AWARD 2144147.

## Notes on contributors

**Xubo Yue** is a PhD candidate in the Department of Industrial & Operations Engineering at the University of Michigan. His research focuses on federated and distributed data analytics. Currently, he is developing federated data analytics methods that rethink how both prescriptive and predictive analytics are achieved within IoT-enabled systems, specifically manufacturing and renewable energy. He has received several best paper awards from the Institute for Operations Research and the Management Sciences (INFORMS), the Institute of Industrial and Systems Engineers (IISE), and other renowned organizations.

**Raed Kontar** is an assistant professor in the Industrial & Operations engineering department at the University of Michigan and an affiliate with the Michigan Institute for Data Science. Raed's research focuses on distributed and federated probabilistic modeling. Raed obtained an undergraduate degree in civil & environmental engineering from the American University of Beirut in 2014, a master's degree in statistics in 2017 and a PhD degree in industrial & systems engineering in 2018, both from the University of Wisconsin-Madison. Raed received the NSF CAREER award in 2022. His research is currently supported by both NSF and NIH.

**Ana María Estrada Gómez** is an assistant professor in the School of Industrial Engineering at Purdue University. She received a BSc in industrial engineering and a B.Sc. in mathematics from la Universidad de los Andes in 2013 and 2015, respectively. She also holds a M.Sc. in industrial engineering from la Universidad de los Andes (2015), and a M.Sc. in statistics from Georgia Tech (2018). In 2021, she received her PhD in industrial engineering with a specialization in statistics from Georgia Tech. Her research interests lie in developing efficient methodologies and algorithms for modeling, monitoring, and diagnosing complex systems collecting high-dimensional data, using statistics and machine learning tools. The methods that she has developed have been applied in the manufacturing, environmental, and healthcare sectors. She is the recipient of the SPES+Q&P Best Student Paper Award from ASA, the QSR Best Poster Award from INFORMS, and the IISE Doctoral Colloquium Best Poster Award. At Georgia Tech, she was recognized with the Graduate Teaching Fellowship, granted by the Center

for Teaching and Learning, and with the Stewart Fellowship, awarded by the School of Industrial and Systems Engineering. She has also been appointed as a Latina Trailblazer in Engineering Fellow by Purdue's College of Engineering.

## References

- Albert, J. and Hu, J. (2019) *Probability and Bayesian Modeling*, CRC Press, Milton Park, Oxfordshire.
- Arashi, M., Roozbeh, M., Hamzah, N.A. and Gasparini, M. (2021) Ridge regression and its applications in genetic studies. *Plos one*, **16**(4), e0245376.
- Barthelme, S. (2016) Simon Barthelme: The expectation-propagation algorithm: A tutorial - part 1. <https://www.youtube.com/watch?v=0tomU1q3AdY&t=2373s>
- Blanco-Filgueira, B., Garcia-Lesta, D., Fernández-Sanjurjo, M., Brea, V.M. and López, P. (2019) Deep learning-based multiple object visual tracking on embedded system for iot and mobile edge computing applications. *IEEE Internet of Things Journal*, **6**(3), 5423–5431.
- CleanTechnica (2021) Tesla fsd hardware has 150 million times more computer power than Apollo 11 computer. <https://cleantechnica.com/2021/05/24/tesla-fsd-hardware-has-150-million-times-more-computer-power-than-apollo-11-computer/#:~:text=Tesla's%20twin%20chips%20and%20combined,and%20Collins%20to%20the%20moon> (accessed 24 May 2021).
- Cortez, P. and Silva, A.M.G. (2008) Using data mining to predict secondary school student performance, in A. Brito and J. Teixeira (eds.), *Proceedings of 5th Annual Future Business Technology Conference*, EUROSIS-ETI, Porto, pp. 5–12.
- Deng, Y., Kamani, M.M. and Mahdavi, M. (2020) Adaptive personalized federated learning. *arXiv preprint arXiv:2003.13461*.
- Du, W., Xu, D., Wu, X. and Tong, H. (2021) Fairness-aware agnostic federated learning. In *Proceedings of the 2021 SIAM International Conference on Data Mining (SDM)*, SIAM, Philadelphia, PA, pp. 181–189.
- Fallah, A., Mokhtari, A. and Ozdaglar, A. (2020) Personalized federated learning: A meta-learning approach. *arXiv preprint arXiv:2002.07948*.
- Fan, J., Guo, Y. and Wang, K. (2021) Communication-efficient accurate statistical estimation. *Journal of the American Statistical Association*, **116**, 1–11.
- Hassan, N., Gillani, S., Ahmed, E., Yaqoob, I. and Imran, M. (2018) The role of edge computing in internet of things. *IEEE Communications Magazine*, **56**(11), 110–115.
- Jordan, M.I., Lee, J.D. and Yang, Y. (2019) Communication-efficient distributed statistical inference. *Journal of the American Statistical Association*, **114**(526), 668–681.
- Kairouz, P., McMahan, H.B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A.N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R. et al. (2021) Advances and open problems in federated learning. *Foundations and Trends in Machine Learning*, **14** (1–2), 1–210.
- Kontar, R., Shi, N., Yue, X., Chung, S., Byon, E., Chowdhury, M., Jin, J., Kontar, W., Masoud, N., Nouiehed, M. et al. (2021) The internet of federated things (ioft). *IEEE Access*, **9**, 156071–156113.
- Lawson, P., McPhail, B. and Lawton, E. (2015) The connected car: Who is in the driver's seat? A study on privacy and onboard vehicle telematics technology. <https://trid.trb.org/view/1348911>
- Li, J., Xu, J. and Zhou, Q. (2018) Monitoring serially dependent categorical processes with ordinal information. *IIE Transactions*, **50**(7), 596–605.
- Li, T., Hu, S., Beirami, A. and Smith, V. (2021) Ditto: Fair and robust federated learning through personalization, in *International Conference on Machine Learning*, Virtual Conference, International Conference on Machine Learning, Vienna, Austria, pp. 6357–6368.
- Li, T., Sahu, A.K., Talwalkar, A. and Smith, V. (2020) Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine*, **37**(3), 50–60.

- Li, T., Sahu, A.K., Zaheer, M., Sanjabi, M., Talwalkar, A. and Smith, V. (2020) Federated optimization in heterogeneous networks. *Proceedings of Machine Learning and Systems*, **2**, 429–450.
- Li, T., Sanjabi, M., Beirami, A. and Smith, V. (2019) Fair resource allocation in federated learning. *arXiv preprint arXiv:1905.10497*.
- Li, X., Huang, K., Yang, W., Wang, S. and Zhang, Z. (2019) On the convergence of fedavg on non-iid data. *arXiv preprint arXiv:1907.02189*.
- Liu, K., Gebraeel, N.Z. and Shi, J. (2013) A data-level fusion model for developing composite health indices for degradation modeling and prognostic analysis. *IEEE Transactions on Automation Science and Engineering*, **10**(3), 652–664.
- Liu, S., Liu, L., Tang, J., Yu, B., Wang, Y. and Shi, W. (2019) Edge computing for autonomous driving: Opportunities and challenges. *Proceedings of the IEEE*, **107**(8), 1697–1716.
- Liu, W., Chen, L., Chen, Y. and Zhang, W. (2020) Accelerating federated learning via momentum gradient descent. *IEEE Transactions on Parallel and Distributed Systems*, **31**(8), 1754–1766.
- McMahan, B., Moore, E., Ramage, D., Hampson, S. and y Arcas, B.A. (2017) Communication-efficient learning of deep networks from decentralized data, in *Artificial Intelligence and Statistics*, The 20th International Conference on Artificial Intelligence and Statistics, Ft. Lauderdale, FL, pp. 1273–1282.
- Minka, T.P. (2001) A family of algorithms for approximate Bayesian inference. Ph.D. thesis, Massachusetts Institute of Technology, Cambridge, MA.
- Minka, T.P. (2013) Expectation propagation for approximate Bayesian inference. *arXiv preprint arXiv:1301.2294*.
- Morell, J.Á. and Alba, E. (2022) Dynamic and adaptive fault-tolerant asynchronous federated learning using volunteer edge devices. *Future Generation Computer Systems*, **133**, 53–67.
- Patterson, D., Gonzalez, J., Le, Q., Liang, C., Munguia, L.-M., Rothchild, D., So, D., Texier, M. and Dean, J. (2021) Carbon emissions and large neural network training. *arXiv preprint arXiv:2104.10350*.
- Rahman, M.A. and Hossain, M.S. (2021) An internet-of-medical-things-enabled edge computing framework for tackling covid-19. *IEEE Internet of Things Journal*, **8**(21), 15847–15854.
- Ruder, S. (2017) An overview of multi-task learning in deep neural networks. *arXiv preprint arXiv:1706.05098*.
- Şahin, D.Ö., Akleyek, S. and Kiliç, E. (2022) Linregdroid: Detection of android malware using multiple linear regression models-based classifiers. *IEEE Access*, **10**, 14246–14259.
- Schulz, M.-A., Yeo, B., Vogelstein, J.T., Mourao-Miranada, J., Kather, J.N., Kording, K., Richards, B. and Bzdok, D. (2020) Different scaling of linear models and deep learning in ukbiobank brain images versus machine-learning datasets. *Nature Communications*, **11**(1), 1–15.
- Shi, N., Lai, F., Kontar, R.A. and Chowdhury, M. (2021) Fed-ensemble: Improving generalization through model ensembling in federated learning. *arXiv preprint arXiv:2107.10663*.
- Shi, W., Cao, J., Zhang, Q., Li, Y. and Xu, L. (2016) Edge computing: Vision and challenges. *IEEE Internet of Things Journal*, **3**(5), 637–646.
- Si, W., Yang, Q. and Wu, X. (2017) A distribution-based functional linear model for reliability analysis of advanced high-strength dual-phase steels by utilizing material microstructure images. *IISE Transactions*, **49**(9), 863–873.
- Song, C. and Liu, K. (2018) Statistical degradation modeling and prognostics of multiple sensor signals via data fusion: A composite health index approach. *IISE Transactions*, **50**(10), 853–867.
- Stich, S.U. (2018) Local SGD converges fast and communicates little. *arXiv preprint arXiv:1805.09767*.
- Tan, A.Z., Yu, H., Cui, L. and Yang, Q. (2022) Towards personalized federated learning. *IEEE Transactions on Neural Networks and Learning Systems*, **33**, 1–17.
- Tibshirani, R. (1996) Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, **58**(1), 267–288.
- Tong, Q., Liang, G., Zhu, T. and Bi, J. (2020) Federated nonconvex sparse learning. *arXiv preprint arXiv:2101.00052*.
- Van Erp, S., Oberski, D.L. and Mulder, J. (2019) Shrinkage priors for Bayesian penalized regression. *Journal of Mathematical Psychology*, **89**, 31–50.
- Vanschoren, J. (2019) Meta-learning, in *Automated Machine Learning*, Springer Nature, Switzerland, pp. 35–61.
- Vehari, A., Gelman, A., Sivula, T., Jylänki, P., Tran, D., Sahai, S., Blomstedt, P., Cunningham, J.P., Schiminovich, D. and Robert, C.P. (2020) Expectation propagation as a way of life: A framework for Bayesian inference on partitioned data. *The Journal of Machine Learning Research*, **21**(17), 1–53.
- Wang, J. and Joshi, G. (2021) Cooperative SGD: A unified framework for the design and analysis of local-update sgd algorithms. *Journal of Machine Learning Research*, **22**(213), 1–50.
- Williams, C.K. and Rasmussen, C.E. (2006) *Gaussian Processes for Machine Learning*, Volume 2. MIT press Cambridge, MA.
- Yu, H., Liu, Z., Liu, Y., Chen, T., Cong, M., Weng, X., Niyato, D. and Yang, Q. (2020) A fairness-aware incentive scheme for federated learning, in *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, Association for Computing Machinery, New York, NY, pp. 393–399.
- Yuan, H. and Ma, T. (2020) Federated accelerated stochastic gradient descent. *Advances in Neural Information Processing Systems*, **33**, 5332–5344.
- Yuan, H., Zaheer, M. and Reddi, S. (2021) Federated composite optimization, in *International Conference on Machine Learning*, Curran Associates Inc., Red Hook, NY, pp. 12253–12266.
- Yue, X. and Kontar, R.A. (2021) Federated Gaussian process: Convergence, automatic personalization and multi-fidelity modeling. *arXiv preprint arXiv:2111.14008*.
- Yue, X., Nouiehed, M. and Kontar, R.A. (2021) Gifair-fl: An approach for group and individual fairness in federated learning. *arXiv preprint arXiv:2108.02741*.
- Zhang, Y., Duchi, J. and Wainwright, M. (2015) Divide and conquer kernel ridge regression: A distributed algorithm with minimax optimal rates. *The Journal of Machine Learning Research*, **16**(1), 3299–3340.
- Zhang, Y. and Yeung, D.-Y. (2012) A convex formulation for learning task relationships in multi-task learning. *arXiv preprint arXiv:1203.3536*.
- Zhu, S., Ota, K. and Dong, M. (2021) Green AI for IIOT: Energy efficient intelligent edge computing for industrial internet of things. *IEEE Transactions on Green Communications and Networking*, **6**, 79–88.