

FedNP: Towards Non-IID Federated Learning via Federated Neural Propagation

Xueyang Wu^{1*}, Hengguan Huang^{2*†}, Youlong Ding³, Hao Wang⁴, Ye Wang², Qian Xu¹

¹Hong Kong University of Science and Technology, Hong Kong SAR, China

²National University of Singapore, Singapore

³Shenzhen University, Shenzhen, China

⁴Rutgers University, Piscataway, NJ, USA

xwuba@connect.ust.hk, huang.hengguan@u.nus.edu, dingyoulon@gmail.com, hw488@cs.rutgers.edu,

wangye@comp.nus.edu.sg, qianxu@ust.hk

Abstract

Traditional federated learning (FL) algorithms, such as FedAvg, fail to handle non-i.i.d data because they learn a global model by simply averaging biased local models that are trained on non-i.i.d data, therefore failing to model the global data distribution. In this paper, we present a novel Bayesian FL algorithm that successfully handles such a non-i.i.d FL setting by enhancing the local training task with an auxiliary task that explicitly estimates the global data distribution. One key challenge in estimating the global data distribution is that the data are partitioned in FL, and therefore the ground-truth global data distribution is inaccessible. To address this challenge, we propose an expectation-propagation-inspired probabilistic neural network, dubbed federated neural propagation (FedNP), which efficiently estimates the global data distribution given non-i.i.d data partitions. Our algorithm is sampling-free and end-to-end differentiable, can be applied with any conventional FL frameworks and learns richer global data representation. Experiments on both image classification tasks with synthetic non-i.i.d image data partitions and real-world non-i.i.d speech recognition tasks demonstrate that our framework effectively alleviates the performance deterioration caused by non-i.i.d data.

Introduction

Federated learning (FL) is an increasingly more important machine learning paradigm where many clients jointly train a powerful global model with cross-silo training data. The major target of FL is to utilize the massive data created and collected by different clients while obeying privacy protection regulations such as the European Union’s General Data Protection Regulation (GDPR) (Voss 2016). The most representative FL algorithm is Federated Averaging (FedAvg) (McMahan et al. 2017a), which successfully trains a powerful global model while keeping training data on each local client (i.e., each mobile phone). Thanks to its feasibility and effectiveness, FedAvg has been successfully adopted in many applications, such as speech recognition (Tan et al. 2021) and language modeling (McMahan et al. 2017b). Recently, more and more FL researchers have

turned their attention to a more practical setting where the data distributions are non-i.i.d (Zhu et al. 2021); in practice, a federation usually consists of different clients from diverse sources whose data distributions are intrinsically distinct. In such non-i.i.d settings, typical FL algorithms often fail to achieve reasonable performance.

Taking a simple polynomial curve fitting task as an example, as shown in Figure 1(a), our goal is to train a model that can fit the data points sampled from a polynomial function. To simulate an extremely non-i.i.d federated data partition, the data points are split into three groups with disjoint ranges of x -coordinates, denoted as $\mathbf{X} = \{\mathbf{X}_k\}_{k=1}^3$. Three federated clients, each training a local model ($\theta = \{\theta_k\}_{k=1}^3$), then collaboratively learn a global model via FL algorithms, e.g., FedAvg. The result shows that while the local models successfully fit their local data in the local training steps, the final prediction is nearly random for test data points. The reason is that the averaging of defective local models adopted by FedAvg cannot provide a global model with sufficient capacity to describe the entire data distribution (see Figure 1(b)). With such insight, we then approach non-i.i.d FL from a new perspective by explicitly modeling a latent global data distribution to enforce local models to have a global view as an auxiliary task.

Existing works in Bayesian federated learning (BFL) (Al-Shedivat et al. 2020; Chen and Chao 2020) focus on directly inferring the global model distribution by aggregating over local model distributions, consequently failing to capture a latent global data distribution. A BFL framework augmented with latent global data distribution should offer several advantages, including (1) a lower dimensional representation of input data uncertainty, caused by limited access to global data; and (2) a greater expressive power to capture the complex dependency underlying the input data across clients than conventional BFL frameworks.

Despite the enormous successes of Bayesian learning and deep learning, it is still challenging to infer such a global data distribution. The reason is that in non-i.i.d federated scenarios, each client only has access to their local data, and therefore the ground-truth global data distribution is inaccessible due to their distinct local data distributions. For instance, Laplace approximation (LA) and Monte-Carlo (MC) approximation, adopted by conventional BFL frameworks (Al-Shedivat et al. 2020; Chen and Chao 2020; Liu

*These authors contributed equally.

†Correspondence to Hengguan Huang

Copyright © 2023, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

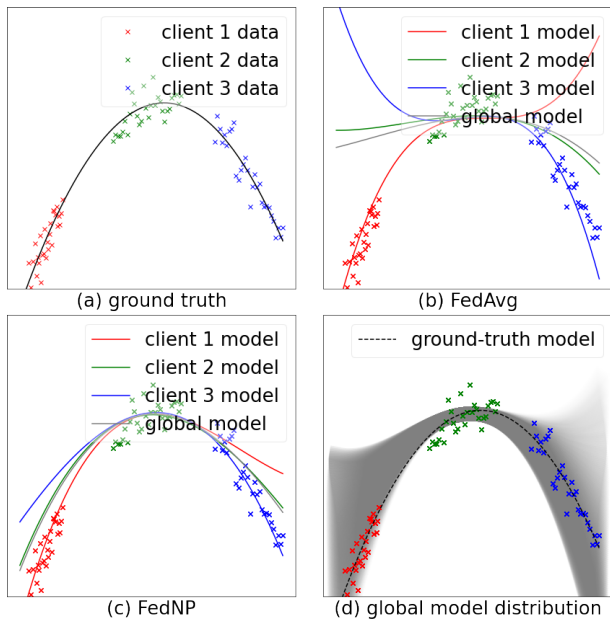


Figure 1: Toy example on polynomial curve fitting task. Data points are denoted by ‘ \times ’ and models are denoted by ‘—’. (a) The ground-truth curve where the observed points are sampled (with Gaussian noise). (b) and (c) The local models of three clients and the global model trained by FedAvg and FedNP, respectively. (d) global model distribution estimated by FedNP (within three standard deviations of the mean).

et al. 2021), are widely adopted techniques to approximate an intractable distribution, but it relies on the observation of global data and therefore cannot apply to our problem settings directly. Another Bayesian inference (BI) method, expectation propagation (EP) (Minka 2013), does not require access to global data and has the potential to approximate the global data distribution. However, applying EP to facilitate our auxiliary task is highly non-trivial:

- EP impractically requires to calculate the intractable likelihood term $\mathcal{L}(\theta_k, \mathbf{X}_k | \mathbf{z})$ due to the unavailability of θ_k , the model parameters for the current mini-batch.
- Since only model parameters for the previous mini-batch are available, we empirically tested an intuitive approach – using such parameters to evaluate the likelihood – and found that the model would not reliably converge.

To overcome these challenges, we reformulate EP by factorizing the latent global data distribution $p(\mathbf{z} | \mathbf{X})$ into several approximate posterior factors $q(\mathbf{z} | \mathbf{X}_k)$, without the need to calculate the likelihood term. We further propose an EP-like probabilistic neural network, dubbed federated neural propagation (FedNP), which follows the general update rule of EP and efficiently estimates the global data distribution. More specifically, FedNP enhances the local training task with an auxiliary task that explicitly estimates a latent global data distribution, stabilizing and correcting the process of local training (See Figure 1(c)), whereby a probabilistic neural network is adopted to map such distribution to a global model distribution, consequently regularizing the

local model by avoiding it sinking into local data distribution. Figure 1(d) shows that FedNP successfully estimates a correct global model parameter distribution and avoids catastrophic failure of local models on unseen data, resulting in a more accurate global model.

Furthermore, unlike existing algorithms for deep-learning-based EP (Jylänki, Nummenmaa, and Vehtari 2014; Soudry, Hubara, and Meir 2014; Heess, Tarlow, and Winn 2013) that require numerical approximation or sampling, we develop a closed-form approximation algorithm for inferring the latent global data distribution, thereby improving the efficiency of modeling federated data.

The major contributions of our work are three-fold:

- We present a new perspective on handling non-i.i.d federated data, which explicitly considers the global data distribution when performing the local training steps.
- We reformulate EP to remove the dependence on the intractable likelihood term and derive a closed-form approximation for the latent global data distribution using deep neural networks, allowing our algorithm to be end-to-end differentiable.
- Our experiments on real-world non-i.i.d image and speech datasets demonstrate that FedNP effectively alleviates the performance deterioration caused by non-i.i.d data compared to state-of-the-art baselines.

Related Work

Federated learning (FL), as a new collaborative learning paradigm, has been gaining more attention in recent years (McMahan et al. 2017a; Yang et al. 2019), and the challenges of FL from non-i.i.d data have been noticed, especially for supervised learning tasks (Li et al. 2021, 2020; Zhao et al. 2018; Xie, Koyejo, and Gupta 2019; Yu et al. 2020). A fundamental idea of these works is to regularize local models during the local training step in FL. Most of them either directly take the averaged model of previous federated turn as the ground truth to regularize the local model training, e.g. FedProx (Li et al. 2020) and MOON (Li, He, and Song 2021), or use a dynamic regularizer, e.g. FedDyn (Acar et al. 2021), FedDC (Gao et al. 2022), or estimate a correction term for local models based on previous aggregated models, e.g. SCAFFOLD (Karimireddy et al. 2020). In contrast, we introduce into local model training an auxiliary task that explicitly estimates the global data distribution from partitioned data, thereby encouraging the local models to be more expressive and preventing them from sinking into local distributions.

From the perspective of Bayesian inference (BI), Bayesian federated learning (BFL) has been studied before, but mainly under the context of model aggregation. For instance, most existing BFL methods (Al-Shedivat et al. 2020; Liu et al. 2021; Chen and Chao 2020) focus on inferring global model distribution by aggregating over local model distributions, while our method aims to regularize local model distributions by inferring a latent global data distribution from input data across clients. Furthermore, similar to FedAvg (McMahan et al. 2017a), FedPA (Al-Shedivat et al. 2020) targets a general federated learning setting, without considering a more challenging extremely

non-i.i.d federated settings, such as ours. To handle non-i.i.d data, FedBE (Chen and Chao 2020) propose a Bayesian ensemble method for server-side model aggregation. However, it requires collecting an unlabeled data set across clients, failing to handle our non-i.i.d FL setting. Other key aspects that distinguish our FedNP from the above methods are as follows. First, FedNP reformulates expectation propagation (EP) with neural networks for handling non-i.i.d FL settings. Inherits from EP, FedNP does not require access to global data during training. In contrast, the Laplace approximation (LA) and Monte-Carlo (MC) approximation, adopted by the above methods, relies on the observation of global data. Second, our FedNP algorithm is sampling-free and end-to-end differentiable, leading to a more efficient and scalable framework. (Please refer to Appendix C for related work on EP with neural networks. *Full appendix and codes can be found at <https://github.com/CodePothon/fednp>.*)

Background

Federated Learning (FL). Following (Li et al. 2019), we formulate FL with the non-i.i.d data distributed on separate clients. We denote the number of clients in the FL system as K and the private annotated dataset in the party $k \in \{1 \dots K\}$ as $(\mathbf{X}_k, \mathbf{Y}_k)$, where $\mathbf{X}_k = \{\mathbf{x}_{k,i}\}_{i=1}^{n_k}$ contains n_k observations and $\mathbf{Y}_k = \{\mathbf{y}_{k,i}\}_{i=1}^{n_k}$ includes the corresponding training labels.

The FL target is to find a global model θ with private data distributed on different clients by the following iterative process (i.e., FedAvg):

1. Each client k optimizes their local model θ_k on its private data set $(\mathbf{X}_k, \mathbf{Y}_k)$ with the objective $J(\cdot)$ according to Eq. 1 below, and send θ_k to the server.
2. The server computes the global model θ_{avg} with weighted average of client models $\{\theta_k\}_{k=1}^K$ with $\theta_{avg} \triangleq \sum_{k=1}^K v_k \theta_k$, where v_k is the weight of the corresponding party k such that $v_k \geq 0$ and $\sum_{k=1}^K v_k = 1$.
3. The server sends θ to all clients and repeat step 1-3 until it reaches the stop criterion.

$$\min_{\theta_k} J(\theta_k; \mathbf{X}_k, \mathbf{Y}_k) \triangleq \frac{1}{n_k} \sum_{i=1}^{n_k} \ell(\mathbf{x}_{k,i}, \mathbf{y}_{k,i}; \theta_k), \quad (1)$$

where $\ell(\cdot; \cdot)$ is the localized loss function, e.g. cross-entropy loss for classification tasks.

Expectation Propagation (EP). Before introducing our FedNP, we begin by reviewing the traditional EP algorithm (Minka 2013). EP is a deterministic approximation algorithm, often used for Bayesian inference of posterior distributions of model parameters, which is believed to be able to provide significantly more accurate approximations than VI (Jordan et al. 1999) and LA (MacKay 1992) methods.

Consider a regression task that predicts some attributes of interest $\theta = \{\theta_k\}_{k=1}^K$ given observations $\mathbf{X} = \{\mathbf{X}_k\}_{k=1}^K$, where K is number of data partitions. Assume both $\{\theta_k\}_{k=1}^K$ and $\{\mathbf{X}_k\}_{k=1}^K$ are conditionally independent given the latent variable \mathbf{z} . As the posterior distribution of interest $p(\mathbf{z}|\mathbf{X})$ is computationally intractable, EP attempts to approximate it with a tractable approximating distribution $q(\mathbf{z})$, which can

be further factorized into multiple approximate factors:

$$q(\mathbf{z}) \propto p_0(\mathbf{z}) \prod_{i=1}^N q_i(\mathbf{z}), \quad (2)$$

where $p_0(\mathbf{z})$ is the prior distribution; the approximate factor $q_i(\mathbf{z})$ is iteratively refined so that they capture the contribution of each data partition \mathbf{X}_k to the posterior $p(\mathbf{z}|\mathbf{X})$:

$$p(\mathbf{z}|\mathbf{X}) \propto p_0(\mathbf{z}) \prod_{k=1}^K \mathcal{L}(\theta_k, \mathbf{X}_k|\mathbf{z}), \quad (3)$$

where $\mathcal{L}(\theta_k, \mathbf{X}_k|\mathbf{z})$ denotes likelihood. If mini-batch optimization is adopted, only mini-batches are used to evaluate such likelihood instead of the entire data partition iteratively. Specifically, EP iterates over the following steps:

1. Construct the cavity distribution by removing one of the approximate factors, i.e., the k -th factor. It can be written as: $q_{-k}(\mathbf{z}) \propto p_0(\mathbf{z}) \prod_{j \neq k} q_j(\mathbf{z})$
2. Integrate the likelihood $\mathcal{L}(\theta_k, \mathbf{X}_k|\mathbf{z})$ to the cavity to produce the hybrid distribution: $h_k(\mathbf{z}) \propto q_{-k}(\mathbf{z}) \mathcal{L}(\theta_k, \mathbf{X}_k|\mathbf{z})$.
3. Update the parameters of the k -th approximated factor $q_k(\mathbf{z})$ through minimizing the KL divergence between the hybrid distribution $h_k(\mathbf{z})$ and the approximated distribution $q(\mathbf{z})$, namely, $\text{KL}[h_k(\mathbf{z}) \| q_{-k}(\mathbf{z}) q_k(\mathbf{z})]$.
4. Update the approximated distribution $q(\mathbf{z})$ by including the updated approximated factor $q_k(\mathbf{z})$: $q(\mathbf{z}) \propto q_{-k}(\mathbf{z}) q_k(\mathbf{z})$.

When $q_k(\mathbf{z})$ is assumed to follow an exponential family distribution (e.g., a Gaussian), minimization of the KL divergence in Step 3 can be reduced to moment matching (Amari and Nagaoka 2000). However, when applying EP to deep neural networks, this moment matching step is computationally intractable, requiring numerical approximation or sampling, and therefore has to compromise between accuracy and efficiency. It's also worth noting that θ_k in our problem settings is the local model parameters for the current mini-batch. Due to its unavailability when evaluating the likelihood term during training, EP cannot be directly applied to our problem settings. Addressing these problems are primary focuses of our model.

Methodology

As shown in the toy experiment, local models tend to sink into local data distributions, resulting in a model drift in the local training step of FL with non-i.i.d data. In the following sections, we firstly formalize such problems and then present a *sampling-free and fully differentiable deep probabilistic neural network* for achieving *end-to-end goals* of inferring latent global data distribution from non-i.i.d partitioned data, thereby regularizing local models and preventing them from sinking into local data distributions.

Problem Formulation

We consider *non-i.i.d* FL settings with classification as major task and aim to enhance the local training steps of FL with an auxiliary task that explicitly models a latent global data distribution to constrain local training. Specifically,

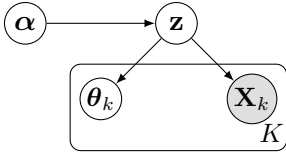


Figure 2: Graphical Model of FedNP. α is the uniform prior of the latent variable \mathbf{z} .

suppose we are given K annotated data partitions $(\mathbf{X}, \mathbf{Y}) = \{(\mathbf{X}_k, \mathbf{Y}_k)\}_{k=1}^K$ privately maintained by K clients, where $\mathbf{X}_k = \{\mathbf{x} | \mathbf{x} \sim p_k(\mathbf{x})\}$ containing data items that are drawn from non-identical distributions, i.e., $p_k(\mathbf{x}) \neq p_{k'}(\mathbf{x})$ for all $k \neq k'$, and \mathbf{Y}_k includes the corresponding training labels. Suppose each client optimizes their local model θ_k on their private data $(\mathbf{X}_k, \mathbf{Y}_k)$. Assume both θ_k and \mathbf{X}_k are conditionally independent given the latent variable \mathbf{z} which represents the global data distribution. The probabilistic model is shown using the standard directed graphical notation in Figure 2. Our model belongs to the broad category of *Bayesian deep learning* (Wang and Yeung 2016, 2020), with a graphical model component as additional inductive bias for the probabilistic neural network. The auxiliary task aims to infer the posterior distribution $p(\mathbf{z} | \mathbf{X})$, whereby the conditional global model distributions $p(\theta_k | \mathbf{z})$ are produced, regularizing each local model θ_k to have a global view.

Federated Neural Propagation

Our proposed Federated Neural Propagation (FedNP) is a special type of message passing algorithm for inferring a latent global data distribution using localized inference over non-i.i.d data partitions and preventing FL local models from sinking into local data distributions. FedNP follows a general updating rule of another message passing algorithm, EP, but adopts a different factorization of the target distribution to remove the dependence on the intractable likelihood. The advantage of FedNP is that it can approximate the global data distribution using fully differentiable neural networks that encode the respective closed-form approximation, rather than inefficient sampling and numerical quadrature adopted by existing neural network-based EP methods.

Specifically, assuming a uniform distribution for the prior $p_0(\mathbf{z})$, the posterior of the latent global data distribution $p(\mathbf{z} | \mathbf{X})$ can be factorized as the product of the posteriors conditioned on local data:

$$p(\mathbf{z} | \mathbf{X}) \propto p_0(\mathbf{z}) \prod_{k=1}^K p(\mathbf{X}_k | \mathbf{z}) \quad (4)$$

$$\propto \prod_{k=1}^K p_0(\mathbf{z}) p(\mathbf{X}_k | \mathbf{z}) \propto \prod_{k=1}^K p(\mathbf{z} | \mathbf{X}_k).$$

Eq. (4) is derived by simply applying Bayes' theorem with the assumption of a uniform prior for $p_0(\theta)$. We then follow EP to adopt K approximate factors to approximate $p(\mathbf{z} | \mathbf{X})$. Specifically, we use $q_k(\mathbf{z})$ as the Gaussian distributed ap-

proximate factor for $p(\mathbf{z} | \mathbf{X}_k)$ and have

$$q(\mathbf{z}) \propto \prod_{k=1}^K q_k(\mathbf{z}), \quad (5)$$

where $q(\mathbf{z})$ approximates global data distribution $p(\mathbf{z} | \mathbf{X})$; $q(\mathbf{z}) = \mathcal{N}(\mathbf{z}_m, \mathbf{z}_s)$, where \mathbf{z}_m and \mathbf{z}_s are the mean and variance, respectively; $q_k(\mathbf{z}) = \mathcal{N}(\mathbf{z}_{m,k}, \mathbf{z}_{s,k})$, where $\mathbf{z}_{m,k}$ and $\mathbf{z}_{s,k}$ are the mean and variance, respectively. Unlike Eq. (2), we omit the prior $p_0(\mathbf{z})$ in the expression as we assume it is a uniform distribution. In the following subsections, we will show how to calculate the mean and variance of $q(\mathbf{z})$ analytically using fully differentiable neural networks.

Infer the Approximate Global Data Distribution $q(\mathbf{z})$. We follow the general updating steps of EP and propose a closed-form approximation of $p(\mathbf{z} | \mathbf{X})$ as $q(\mathbf{z})$. We firstly initialize mean and variance of each Gaussian factor, $q_k(\mathbf{z}) = \mathcal{N}(\mathbf{z}_{m,k}, \mathbf{z}_{s,k})$. We follow Step 1 of EP in Section **Background** to construct the cavity and the hybrid distributions $q_{-k}(\mathbf{z})$, $h_k(\mathbf{z})$:

$$q_{-k}(\mathbf{z}) \propto \prod_{j \neq k} q_j(\mathbf{z}), \quad (6)$$

where the cavity distribution $q_{-k}(\mathbf{z}) = \mathcal{N}(\mathbf{z}_{m,-k}, \mathbf{z}_{s,-k})$, whose mean $\mathbf{z}_{m,-k}$ and variance $\mathbf{z}_{s,-k}$ are calculated by taking product of Gaussian factors.

We then follow Step 2 of EP in Section **EP** to construct the hybrid distribution $h_k(\mathbf{z})$:

$$h_k(\mathbf{z}) \propto p(\mathbf{z} | \mathbf{X}_k) q_{-k}(\mathbf{z}), \quad (7)$$

Inspired by nature-parameter-network (NPN) (Wang, Shi, and Yeung 2016), to parameterize the posterior $p(\mathbf{z} | \mathbf{X}_k)$, we adopt deep neural networks to probabilistically propagate information from \mathbf{X}_k to the latent variable \mathbf{z} :

$$p(\mathbf{z} | \mathbf{X}_k) = \phi(\mathbf{z}, \mathbf{X}_k), \quad (8)$$

where ϕ is a neural network (see implementation details in Appendix of corresponding experiment settings).

When following Step 3 and 4 of EP in Section **EP** to update $q(\mathbf{z})$, we find that moment matching is not analytical, due to deep neural networks involved in $\phi(\mathbf{z}, \mathbf{X}_k)$. To overcome this challenge, we propose the following theorem which provides a closed-form solution for moment matching, and thus obtaining closed-form approximation, $q(\mathbf{z})$, for the posterior of the global data distribution, $p(\mathbf{z} | \mathbf{X})$.

Theorem 1. *Suppose we are given a data partition \mathbf{X}_k located at the k -th party during FL. Assume data partitions $\{\mathbf{X}_k\}_{k=1}^K$ are conditionally independent given a latent variable \mathbf{z} . Let $f : \mathcal{R}^{|\mathbf{X}_k|} \rightarrow \mathcal{R}^1$ be a neural network taking as input \mathbf{X}_k . Let $C = f(\mathbf{X}_k)\mathbf{z}$ and $C \sim \mathcal{N}(C_{m,-k}, C_{s,-k})$. Let $q_{-k}(\mathbf{z})$ and $h_k(\mathbf{z})$ be the cavity distribution (defined in Eq. (6)) and the hybrid distribution (defined in Eq. (7)), respectively. We further define $q(\mathbf{z}) = \mathcal{N}(\mathbf{z}_m, \mathbf{z}_s)$ as in Eq. (5). There exists a function $\phi : \mathcal{R}^{|\mathbf{z}|} \times \mathcal{R}^{|\mathbf{X}_k|} \rightarrow \mathcal{R}^1$, such that the update rules of \mathbf{z}_m and \mathbf{z}_s can be written in closed form as:*

$$\mathbf{z}_m = S_1, \quad (9)$$

$$\mathbf{z}_s = S_2 - S_1^2, \quad (10)$$

where

$$S_1 = [(C_{m,-k} + C_{s,-k})E_C(\sigma(C)) - C_{s,-k}E_C(\sigma^2(C))] / S_0 f(\mathbf{X}_k), \quad (11)$$

$$S_2 = [(C_{m,-k} + 2C_{s,-k})E_C(\sigma^2(C)) - 2C_{s,-k}E_C(\sigma^3(C))] / S_0 f^2(\mathbf{X}_k), \quad (12)$$

$$S_0 = E_C(\sigma(C)). \quad (13)$$

We leave the proof in Appendix A.2.

Given a data partition \mathbf{X}_k , Theorem 1 provides the closed-form updates for mean and variance of $q(\mathbf{z}) = \mathcal{N}(\mathbf{z}_m, \mathbf{z}_s)$. Note that Theorem 1 requires an auxiliary distribution of the latent variable $C = f(\mathbf{X}_k)\mathbf{z}$. If the cavity distribution is Gaussian, i.e., $q_{-k}(\mathbf{z}) = \mathcal{N}(\mathbf{z}_{m,-k}, \mathbf{z}_{s,-k})$, C also follows a Gaussian distribution. Therefore denoting this distribution of C as $q_{c,-k}(C)$, we have that:

$$q_{c,-k}(C) = \mathcal{N}(C_{m,-k}, C_{s,-k}) \times \prod_{j \neq k} q_j(C), \quad (14)$$

where $\mathcal{N}(C_{m,-k}, C_{s,-k})$ is a Gaussian distribution with the mean and variance:

$$[C_{m,-k}, C_{s,-k}]^\top = [f(\mathbf{X}_k)\mathbf{z}_{m,-k}, f^2(\mathbf{X}_k)\mathbf{z}_{s,-k}]^\top. \quad (15)$$

Theorem 1 further requires the first three moments of $q_{c,-k}(C)$. We therefore adopt Theorem 2 below to calculate these first three moments, i.e., $E_C(\sigma(C))$, $E_C(\sigma^2(C))$, and $E_C(\sigma^3(C))$, thereby providing an analytical solution for calculating the first two moments of the hybrid distribution $h_k(\mathbf{z})$, S_1 and S_2 . Theorem 1 then follows Step 4 (Section EP) to update $q(\mathbf{z}) = \mathcal{N}(\mathbf{z}_m, \mathbf{z}_s)$, whose parameters \mathbf{z}_m and \mathbf{z}_s can be updated in closed form as well.

Theorem 2. *Suppose $C \sim \mathcal{N}(C_{m,-k}, C_{s,-k})$. Let $d \geq 1$ be a positive integer. There exist two real constants a and b , such that the first d moments can be expressed in closed form:*

$$E_C(\sigma^d(C)) \approx \sigma \left(\frac{a(C_{m,-k} + b)}{\sqrt{1 + \zeta^2 a^2 C_{s,-k}}} \right). \quad (16)$$

The proof for Theorem 2 is in Appendix B.

Generally, Theorem 1 & 2 allow our FedNP algorithm to be end-to-end differentiable. We present such an algorithm with a mild extension to the commonly used FL framework (FedAvg) in Algorithm 1 (see Appendix B).

Infer the Conditional Global Model Distribution $p(\theta_k|\mathbf{z})$. In this section, we aim to infer the conditional global model distribution $p(\theta_k|\mathbf{z})$ based on the approximate global data distribution $q(\mathbf{z})$. Since natural parameter network (NPN) (Wang, Shi, and Yeung 2016) allows inferring a target distribution based on an input distribution through layers of efficient sampling-free probabilistic transformation, we adopt such a model to infer $p(\theta_k|\mathbf{z})$ given $q(\mathbf{z})$, i.e.:

$$p(\theta_k|\mathbf{z}) = \text{NPN}(q(\mathbf{z})), \quad (17)$$

where $\text{NPN}(\cdot)$ represents the natural parameter network, whose implementation is detailed in Appendix.

Learning

Our design of the loss function aims at regularizing local models to have a global view using the conditional global

model distribution $p(\theta_k|\mathbf{z})$. Since we only have θ_k for the previous mini-batch, we take the conditional global model distribution $p(\theta_k|\mathbf{z})$ to match an auxiliary Gaussian distribution $p(\hat{\theta}_k) = \mathcal{N}(\theta_k, \epsilon)$, where ϵ is a vector with all entries set as small constants. Following (NPN) (Wang, Shi, and Yeung 2016), the loss can be written as:

$$\ell = KL [p(\theta_k|\mathbf{z}) \| p(\hat{\theta}_k)] \quad (18)$$

To jointly train FedNP and local models of FL, our final loss for local training can be written as:

$$J(\theta_k; \mathbf{X}_k, \mathbf{Y}_k) + \lambda \ell, \quad (19)$$

where $J(\cdot)$ is the loss for local training (Eq. 1) and λ is the hyperparameter balancing two losses.

Approximation Error and Computational Efficiency. In Theorem 2, we use the probit function $\Phi(\zeta a(C + b))$ to approximate $\sigma^d(C)$. Similar approximation is also adopted by (Wang, Shi, and Yeung 2016) and (Wang and Manning 2013). Both the numerical and theoretical approximation analysis has been well studied in Section 3.2 (Wang, Shi, and Yeung 2016). Our closed-form update steps are efficient in computation. Suppose that there is a system with K data partitions/clients (note that K is usually known beforehand) and the likelihood term is tractable, for EP with neural networks in (Jylänki, Nummenmaa, and Vehtari 2014) which requires numerical quadratures, its computational complexity is $O(MK)$, where M is the number of quadratures points. For EP with MCMC, its computational complexity is $O(NK)$, where N is the number of MCMC samples. To approach a good approximation, both N and M should be sufficiently large (Barthelmé and Chopin 2014). In contrast, with our closed-form update, FedNP’s computational complexity is reduced to $O(K)$.

We show the computational efficiency as well as the approximation accuracy in Appendix D.3.

Experiments

In this section, we evaluate our FedNP on non-i.i.d. cross-silo datasets on numerical regression, image classification, and speech recognition tasks to demonstrate its effectiveness. We describe the experimental details, including the environments, implementations, etc., in Appendix D. The datasets used in our experiments are public, and codes can be found in the supplementary file.

Toy Experiments

Dataset. We evaluate our FedNP on a toy cross-silo non-i.i.d dataset. Given x , we will use the following quadratic polynomial function to generate the labels.

$$y = -x^2 + 2x + 5 + \epsilon \quad (20)$$

where $\epsilon \sim \mathcal{N}(0, 5)$. We assume three clients and sample data points from three disjoint segments as the local training data for each client (see Figure 1(a)).

Results. Figure 1(b) visualizes the trained local models and global model of FedAvg. The local training leads the model to learn a biased curve due to the skewed local data distribution. Although all the three local models fit their own data points perfectly, they fail to generalize well to the unseen global data. This is because the global model obtained

by averaging fails to capture the global data distribution. In contrast, our FedNP is more robust (see Figure 1(c)), and the estimated global model distribution is more accurate (Figure 1(d)), which corrects the local training and avoids performance deterioration. The experiment setup, quantitative results along with qualitative comparison with FedPA are shown in Appendix D.2.

Image Classification with Non-IID Image Datasets

Task. We evaluate the performance of FedNP on image classification, which is the most fundamental task in image datasets under the supervised learning setting. We focus on a setting where the local data distributions across parties are non-i.i.d, which increases the difficulty of training.

Datasets. We conduct experiments on CIFAR-100 (Krizhevsky, Nair, and Hinton 1995) and TinyImageNet (Le and Yang 2015). Similar to previous study (Wang et al. 2020), we use Dirichlet distribution to generate the non-i.i.d data partition among clients. With the above partitioning strategy, each client may have relatively few data samples in some classes. We leave details in Appendix D.3.

Baselines. We compare FedNP with four approaches including (1) FedAvg (McMahan et al. 2017a), (2) FedProx (Li et al. 2020), (3) MOON (Li, He, and Song 2021), (4) SCAFFOLD (Karimireddy et al. 2020), (5) FedDyn (Acar et al. 2021), (6) FedDC (Gao et al. 2022), (7) FedPA (Al-Shedivat et al. 2020), and (8) FedLA (Liu et al. 2021) as discussed in related work. For more details, please refer to Appendix D.3. To further validate the robustness and scalability of FedNP, we vary the number of clients from 10, 50, and 100. To keep the main paper concise, we leave the full performance table in Appendix D.3 and show the performance on the 10-client default setting in Table 1.

Results. Table 1 shows the top-1 test accuracy of all approaches with the 10-client default setting. Comparing different FL approaches, we can observe that FedNP is consistently the best approach. It outperforms FedAvg by around 2% accuracy on average. FedProx’s accuracy is worse than FedAvg. This is because directly minimizing the distance between the global model and the local model will negatively affect convergence (note that the initial local model in the local training phase is exactly the server-side model). MOON also minimizes the distance between the global model and local models. Its difference from FedProx is that MOON defines the distance in the feature space instead of in the parameter space. FedDyn’s and FedDC’s dynamic regularizers do not offer significant performance gain for sophisticated models. Therefore, their performances are also close to FedAvg. The theoretical guarantee of SCAFFOLD relies on the strong smoothness assumption, which does not necessarily hold true in deep learning. Therefore, its performance suffers a lot on these image classification tasks, which has also been verified by Li et al. (2021). Compared with the state-of-the-art BFL frameworks, our FedNP outperforms FedLA (Liu et al. 2021) and FedPA (Al-Shedivat et al. 2020), which employ Laplace approximation/MC-approximation for model aggregation, by around 5% and 2% in terms of accuracy, respectively. This demonstrates the im-

portance of inferring a latent global data distribution with a sampling-free and end-to-end differentiable BFL framework. Generally, our experimental results on CIFAR 100 and TinyImageNet show superior performance and demonstrate the effectiveness of the proposed FedNP.

Methods	CIFAR100 \uparrow	TinyImageNet \uparrow
FedAvg	62.93% \pm 0.3%	51.89% \pm 0.5%
FedProx	62.11% \pm 0.2%	50.63% \pm 0.2%
MOON	63.07% \pm 0.3%	51.44% \pm 0.3%
SCAFFOLD	53.50% \pm 0.3%	46.28% \pm 0.2%
FedDyn	62.04% \pm 0.3%	47.34% \pm 0.2%
FedDC	63.15% \pm 0.3%	47.34% \pm 0.2%
FedPA	63.44% \pm 0.3%	49.57% \pm 0.2%
FedLA	60.98% \pm 0.4%	50.23% \pm 0.2%
FedNP (ours)	65.03% \pm 0.2%	53.18% \pm 0.3%

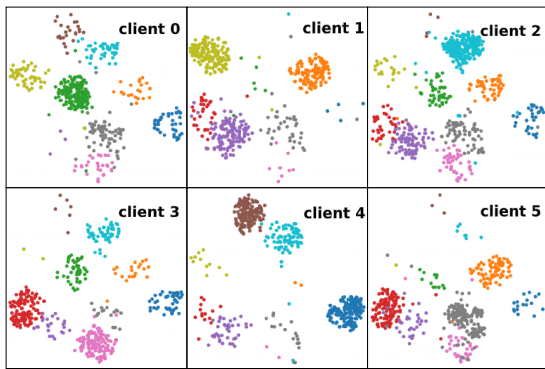
Table 1: The top-1 accuracy of FedNP and the other baselines of the 10-client setting on test datasets. We run three trials and report the mean and standard deviation.

To facilitate the qualitative evaluation of our FedNP, we randomly sample data from the first 10 classes of CIFAR-100 and utilize t-SNE (Van der Maaten and Hinton 2008) to visualize the corresponding ResNet18 backbone (served as the shared encoder for classifier and $q(\mathbf{z})$ output produced by FedNP in Figure 3. Figure 3 shows two advantages of our FedNP compared to FedAvg. First, FedNP is capable of learning more discriminative features among classes with more inner-class compactness and larger inter-class margins. Second, FedNP is able to preserve the structure of global data distribution even in the clients (e.g., client 1 in Figure 3(b)) whose class distributions are extremely imbalanced; in contrast, the data distribution learned by FedAvg is relatively messy. Besides, we visualize the predicted model distribution by FedNP in Figure 5 to verify the efficacy of estimating the global model distribution, as shown in Appendix D.3. Moreover, to evaluate the computational efficiency and accuracy of the approximation for the proposed closed-form update, we compare FedNP with the closed-form update and its sampling-based variant. We leave the details and results in Appendix D.3.

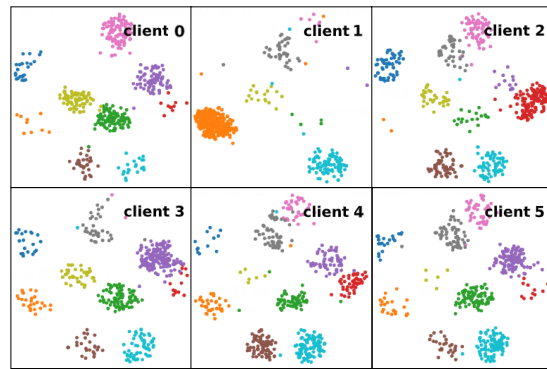
Speech Recognition with Large-scale Non-IID Conversational Speech Dataset

Task. We evaluate our method on a speech task to demonstrate the efficacy of our proposed FedNP on natural real-world non-i.i.d speech datasets. We applied FedNP to the learning of deep neural network acoustic models. The goal of this task is to transcribe a piece of speech to text.

Datasets. We evaluate our proposed method on a challenging real conversational speech dataset CHiME-5 (Barker et al. 2018), whose data are collected from daily life with diverse environments and various speakers. CHiME-5 is a large-scale corpus of real-world multi-speaker conversational speech in home environments. The training set is originally collected by 16 conversation sessions, each consisting



(a) Features of local images extracted by the FedAvg model.



(b) Features of local images extracted by the FedNP model.

Figure 3: T-SNE visualizations on features of local images extracted by the models trained via FedAvg and FedNP on CIFAR-100. For better presentation, only the first 10 classes (with the class id of 0, 1, ..., 9) and the first 6 clients are presented.

of different speakers at different places and talking about different topics; these sessions compose a natural non-i.i.d. data partition. Hence, the non-i.i.d. federated speech setting has 16 clients, one for each session during training, and applies the original testing set, containing 4 conversation sessions. The detailed composition of 16 sessions is reported in Appendix D.4.

Baselines. The baseline speech recognition systems are deep neural network acoustic models on clients trained using typical FL algorithms. Specifically, the deep neural networks are trained with an SRU-HMM-based acoustic model scheme (Huang et al. 2020, 2021), where SRU is a popular and efficient recurrent neural network for acoustic modeling.

Methods	WER ↓
FedAvg (McMahan et al. 2017a)	68.86±0.11
FedProx (Li et al. 2020)	66.72±0.28
FedPA (Al-Shedivat et al. 2020)	69.92±0.22
FedNP (ours)	64.30±0.21

Table 2: The WERs of different methods on CHiME-5 under non-i.i.d. FL settings. The WERs are reported with the mean and standard deviation of three trials.

In the implementation, SRU is used as a backbone of the acoustic model, which contains 12 stacked layers with 1280 hidden nodes. We compare FedNP with three baselines: (1) FedAvg (McMahan et al. 2017a), (2) FedProx (Li et al. 2020), and (3) FedPA (Al-Shedivat et al. 2020). SCAFFOLD and FedLA are unstable in our preliminary experiments due to the large models and complex tasks; therefore, we do not include them as baselines. All baselines and our FedNP adopt identical federated configurations. More implementation details can be found in Appendix D.4.

Results. We train the models using three different random seeds and report the mean and standard deviation (STD) of the test word error rates (WERs) in Table 2. Our FedNP achieves around 4.56 % absolute WER reduction compared

to the FedAvg baseline, a large margin in ASR, suggesting that our FedNP is capable of improving the performance under realistic non-i.i.d. speech datasets. Notably, FedPA performs slightly worse than FedAvg due to its instability. Even compared to the state-of-the-art non-i.i.d. FL algorithm FedProx, our FedNP achieves around 2.42 % absolute WER reduction, which demonstrates that compared to naively typical methods that naively push local models closer to the global model averaged in the last turn, FedNP can estimate a more accurate global model distribution, leading to better performance for non-i.i.d. data.

Conclusion

We present a novel sampling-free and end-to-end differentiable Bayesian federated learning framework, dubbed FedNP, for non-i.i.d. cross-silo data by enhancing the local model with an auxiliary task that explicitly estimates the global data distribution. We successfully alleviate the critical challenge in estimating the global data distribution on partitioned non-i.i.d. data with an expectation-propagation-inspired probabilistic neural network. More specifically, we tackle defects of existing algorithms for deep-learning-based EP inference and derive a closed-form solution for estimating the global data distribution, leading to a more efficient solution for global data modeling. Experiments on toy non-i.i.d. data and real-world extremely non-i.i.d. image and speech data partitions demonstrate that our framework effectively alleviates the performance deterioration caused by non-i.i.d. data compared to other representative baselines.

So far, FedNP has not been evaluated on the unsupervised learning setting, nor on datasets with other modalities such as text. Another current limitation of FedNP is that though FedNP helps facilitate the joint modeling across data centers as an FL algorithm, such cooperation might also bring the risk of personal privacy leakage. Therefore, it is also necessary to explore more privacy-protection methods that can further improve the data privacy and security protection of FL algorithms, including FedNP.

Acknowledgements

The authors thank the reviewers/SPC/AC for the constructive comments to improve the paper. HH and YW are partially supported by a grant A-0008153-00-00 from Ministry of Education of Singapore. HW is partially supported by NSF Grant IIS-2127918 and an Amazon Faculty Research Award. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of the sponsors.

References

- Acar, D. A. E.; Zhao, Y.; Navarro, R. M.; Mattina, M.; Whatmough, P. N.; and Saligrama, V. 2021. Federated learning based on dynamic regularization. *arXiv preprint arXiv:2111.04263*.
- Al-Shedivat, M.; Gillenwater, J.; Xing, E.; and Rostamizadeh, A. 2020. Federated learning via posterior averaging: A new perspective and practical algorithms. *arXiv preprint arXiv:2010.05273*.
- Amari, S.-i.; and Nagaoka, H. 2000. *Methods of information geometry*, volume 191. American Mathematical Soc.
- Barker, J.; Watanabe, S.; Vincent, E.; and Trmal, J. 2018. The fifth 'CHiME' speech separation and recognition challenge: dataset, task and baselines. *arXiv preprint arXiv:1803.10609*.
- Barthelmé, S.; and Chopin, N. 2014. Expectation propagation for likelihood-free inference. *Journal of the American Statistical Association*, 109(505): 315–333.
- Chen, H.-Y.; and Chao, W.-L. 2020. Fedbe: Making bayesian model ensemble applicable to federated learning. *arXiv preprint arXiv:2009.01974*.
- Gao, L.; Fu, H.; Li, L.; Chen, Y.; Xu, M.; and Xu, C.-Z. 2022. FedDC: Federated Learning with Non-IID Data via Local Drift Decoupling and Correction. In *IEEE Conference on Computer Vision and Pattern Recognition*.
- Heess, N.; Tarlow, D.; and Winn, J. 2013. Learning to pass expectation propagation messages. *Advances in Neural Information Processing Systems*, 26: 3219–3227.
- Huang, H.; Liu, H.; Wang, H.; Xiao, C.; and Wang, Y. 2021. STRODE: Stochastic Boundary Ordinary Differential Equation. In *International Conference on Machine Learning*, 4435–4445. PMLR.
- Huang, H.; Xue, F.; Wang, H.; and Wang, Y. 2020. Deep graph random process for relational-thinking-based speech recognition. In *International Conference on Machine Learning*, 4531–4541. PMLR.
- Jordan, M. I.; Ghahramani, Z.; Jaakkola, T. S.; and Saul, L. K. 1999. An introduction to variational methods for graphical models. *Machine learning*, 37(2): 183–233.
- Jylänki, P.; Nummenmaa, A.; and Vehtari, A. 2014. Expectation Propagation for Neural Networks with Sparsity-Promoting Priors. *Journal of Machine Learning Research*, 15(54): 1849–1901.
- Karimireddy, S. P.; Kale, S.; Mohri, M.; Reddi, S.; Stich, S.; and Suresh, A. T. 2020. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, 5132–5143. PMLR.
- Krizhevsky, A.; Nair, V.; and Hinton, G. 1995. CIFAR-100 (Canadian Institute for Advanced Research). *Canadian Institute for Advanced Research*.
- Le, Y.; and Yang, X. 2015. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7): 3.
- Li, Q.; Diao, Y.; Chen, Q.; and He, B. 2021. Federated learning on non-iid data silos: An experimental study. *arXiv preprint arXiv:2102.02079*.
- Li, Q.; He, B.; and Song, D. 2021. Model-Contrastive Federated Learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10713–10722.
- 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, 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, L.; Zheng, F.; Chen, H.; Qi, G.-J.; Huang, H.; and Shao, L. 2021. A Bayesian Federated Learning Framework with Online Laplace Approximation. *arXiv preprint arXiv:2102.01936*.
- MacKay, D. J. 1992. A practical Bayesian framework for backpropagation networks. *Neural computation*, 4(3): 448–472.
- McMahan, B.; Moore, E.; Ramage, D.; Hampson, S.; and y Arcas, B. A. 2017a. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, 1273–1282. PMLR.
- McMahan, H. B.; Ramage, D.; Talwar, K.; and Zhang, L. 2017b. Learning differentially private recurrent language models. *arXiv preprint arXiv:1710.06963*.
- Minka, T. P. 2013. Expectation propagation for approximate Bayesian inference. *arXiv preprint arXiv:1301.2294*.
- Soudry, D.; Hubara, I.; and Meir, R. 2014. Expectation backpropagation: Parameter-free training of multilayer neural networks with continuous or discrete weights. In *NIPS*, volume 1, 2.
- Tan, C.; Jiang, D.; Peng, J.; Wu, X.; Xu, Q.; and Yang, Q. 2021. A de novo divide-and-merge paradigm for acoustic model optimization in automatic speech recognition. In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, 3709–3715.
- Van der Maaten, L.; and Hinton, G. 2008. Visualizing data using t-SNE. *Journal of machine learning research*, 9(11).
- Voss, W. G. 2016. European union data privacy law reform: General data protection regulation, privacy shield, and the right to delisting. *The Business Lawyer*, 72(1): 221–234.
- Wang, H.; Shi, X.; and Yeung, D.-Y. 2016. Natural-parameter networks: A class of probabilistic neural networks. *arXiv preprint arXiv:1611.00448*.

- Wang, H.; and Yeung, D.-Y. 2016. Towards Bayesian deep learning: A framework and some existing methods. *TDKE*, 28(12): 3395–3408.
- Wang, H.; and Yeung, D.-Y. 2020. A Survey on Bayesian Deep Learning. *CSUR*, 53(5): 1–37.
- Wang, H.; Yurochkin, M.; Sun, Y.; Papailiopoulos, D.; and Khazaeni, Y. 2020. Federated learning with matched averaging. *arXiv preprint arXiv:2002.06440*.
- Wang, S.; and Manning, C. 2013. Fast dropout training. In *international conference on machine learning*, 118–126. PMLR.
- Xie, C.; Koyejo, S.; and Gupta, I. 2019. Asynchronous federated optimization. *arXiv preprint arXiv:1903.03934*.
- Yang, Q.; Liu, Y.; Chen, T.; and Tong, Y. 2019. Federated Machine Learning: Concept and Applications. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 10(2): 12.
- Yu, F.; Rawat, A. S.; Menon, A.; and Kumar, S. 2020. Federated learning with only positive labels. In *International Conference on Machine Learning*, 10946–10956. PMLR.
- Zhao, Y.; Li, M.; Lai, L.; Suda, N.; Civan, D.; and Chandra, V. 2018. Federated learning with non-iid data. *arXiv preprint arXiv:1806.00582*.
- Zhu, H.; Xu, J.; Liu, S.; and Jin, Y. 2021. Federated Learning on Non-IID Data: A Survey. 1–29.