

Calibrated One Round Federated Learning with Bayesian Inference in the Predictive Space

Mohsin Hasan^{1,2}, Guojun Zhang³, Kaiyang Guo³, Xi Chen³, Pascal Poupart^{1,2}

¹University of Waterloo

²Vector Institute

³Huawei Noah's Ark Lab

mohsin.hasan@uwaterloo.ca, {guojun.zhang, guo.kaiyang, xi.chen4}@huawei.com, ppoupart@uwaterloo.ca

Abstract

Federated Learning (FL) involves training a model over a dataset distributed among clients, with the constraint that each client's dataset is localized and possibly heterogeneous. In FL, small and noisy datasets are common, highlighting the need for well-calibrated models that represent the uncertainty of predictions. The closest FL techniques to achieving such goals are the Bayesian FL methods which collect parameter samples from local posteriors, and aggregate them to approximate the global posterior. To improve scalability for larger models, one common Bayesian approach is to approximate the global predictive posterior by multiplying local predictive posteriors. In this work, we demonstrate that this method gives *systematically overconfident* predictions, and we remedy this by proposing β -Predictive Bayes, a Bayesian FL algorithm that interpolates between a mixture and product of the predictive posteriors, using a tunable parameter β . This parameter is tuned to improve the global ensemble's calibration, before it is distilled to a single model. Our method is evaluated on a variety of regression and classification datasets to demonstrate its superiority in calibration to other baselines, even as data heterogeneity increases. Code available at <https://github.com/hasanmohsin/betaPredBayesFL>.

Introduction

Federated learning (FL) is a machine learning paradigm that trains a statistical model using decentralized data stored on client devices, with the constraint that client data is kept local (McMahan et al. 2017). FL has found use in smartphone applications, as well as fields such as healthcare and finance due to the abundance of use-cases where sensitive training data is owned by separate entities (Zheng et al. 2022).

The workflow of a typical FL algorithm involves the local training of models on each client, followed by the communication and aggregation of these into a single global model on a central server. Many FL techniques alternate between these two steps until some notion of convergence is reached.

For the purpose of this work, we are concerned with the design of FL algorithms with three key metrics in mind:

1. **Communication Cost:** the transmission of models between clients and servers can be expensive, especially when each model has numerous parameters (e.g. neural

networks). Cost effective FL techniques therefore aim to maximize the amount of local computation at each client, while reducing the rounds of communication.

2. **Performance with Heterogeneous Data:** clients may have different data distributions from each other. This causes the local client models to drift apart, presenting issues when aggregating them into a global model (Zhao et al. 2018). In particular, when client datasets are heavily heterogeneous, global models tend to perform poorly on clients' local datasets.
3. **Calibration:** Clients' data may have too few training points, or too much variance. Therefore, a valuable goal in FL is to produce models that are **well-calibrated**, or in other words: make probabilistic predictions with accurate uncertainty estimates.

Many techniques frame FL as a distributed optimization problem. For these, there is typically a trade-off between achieving a low communication cost and good performance on heterogeneous data: they improve global performance at the cost of more communication (McMahan et al. 2017; Woodworth, Patel, and Srebro 2020). Moreover, they do not have any systematic way of calibrating the global model.

To alleviate these problems, an alternative branch of FL techniques takes a Bayesian perspective on learning the model. They approximate the Bayesian posterior distribution of each client model, and aggregate them into the global Bayesian posterior (Al-Shedivat et al. 2021; Neiswanger, Wang, and Xing 2014). This allows the training of more effective global models on heterogeneous data by leveraging client uncertainty during aggregation (Al-Shedivat et al. 2021). Certain Bayesian FL methods also demonstrate how to aggregate local posteriors in only a single round of communication (Neiswanger, Wang, and Xing 2014). Furthermore, these methods explicitly represent and adjust the uncertainty in model parameters, i.e., they should be well-calibrated in principle.

In these methods, the Bayesian posterior is a distribution over model parameters, and it is expensive to represent and manipulate. Therefore, it becomes necessary to apply approximations to the posterior. Many methods work with crude approximations to the client posterior (as e.g., a multivariate Gaussian) (Al-Shedivat et al. 2021; Neiswanger, Wang, and Xing 2014), which can incur heavy error, result-

ing in poor calibration and accuracy.

A promising Bayesian method for FL is the Bayesian Committee Machine (BCM, Tresp 2000), which instead aggregates the distribution over *model predictions* (referred to as the **Bayesian predictive posterior**) rather than the posterior over *parameters*. The former allows for more accurate approximation due to its lower dimensionality. Nevertheless, these methods need to rely on an approximate aggregation technique, which can add bias to the global model.

In this work, we argue that the BCM is poorly calibrated due to producing overconfident predictions. We remedy this by proposing a new aggregation method that interpolates between the BCM predictions, and those made by a mixture model over the local predictive posteriors. This server aggregation results in an ensemble model, which is then distilled into a single model, to be sent back to clients.

The primary contributions of our work are as follows:

- We propose a novel algorithm for Bayesian FL, called *β -Predictive Bayes*. This method operates in a *single round of communication*, while benefiting from performance over heterogeneous data like the BCM. On the other hand, it does not suffer from the same calibration issues as the BCM.
- We empirically evaluate *β -Predictive Bayes* on multiple regression and classification datasets, using partitions simulating varying levels of data heterogeneity. The proposed technique competes with or outperforms other baselines with respect to calibration, particularly when data heterogeneity increases.

Background

In federated learning (FL), data is distributed across several clients. Let $\mathcal{D} = \mathcal{D}_1 \cup \dots \cup \mathcal{D}_n$ where $\mathcal{D}_i = \{(x_1, y_1), \dots, (x_{k_i}, y_{k_i})\}$ is the dataset of size k_i at client i , consisting of inputs x and targets y . The goal is to learn a predictive model without any data leaving each client to preserve privacy. Let m_θ denote a model, parameterized by weights θ , which outputs a predictive distribution $m_\theta(x) = p(y|x, \theta)$. In a typical FL algorithm, each client learns a local model m_{θ_i} , which it shares with a trusted server that aggregates them into a global model $m_{\bar{\theta}}$. For example, Federated Averaging (FedAvg) (McMahan et al. 2017) aggregates local models by taking the average of their parameters (i.e., $\bar{\theta} = \sum_i \theta_i k_i / (\sum_i k_i)$).

In practice, the datasets are often heterogeneous, which means that they are sampled from different distributions. To deal with heterogeneity and avoid client divergence, FedAvg (McMahan et al. 2017) and many other variants (Mohri, Sivek, and Suresh 2019; Wang et al. 2019; Li et al. 2020a; Wang et al. 2020; Li et al. 2020b) perform frequent rounds of model updates and averaging where in each round the clients update their local models based on a few steps of gradient descent, or more generally, some form of partial training. This can be quite costly due to the increased communication and the need for synchronization at each round. Some methods alter the training to improve performance in heterogeneous settings. For instance, FedProx (Li et al. 2020a) adds a penalty between the global model

and each client’s local model to the local training loss, while Reddi et al. (2021) use an adaptive optimizer (such as Adam (Kingma and Ba 2015)) at the server for aggregation.

Bayesian Learning

Bayesian learning is a training method that takes into account uncertainty over parameters (Bishop 2006). It operates by setting a “model space prior” $p(\theta)$, which encapsulates beliefs about the model parameters before observing data. Upon processing the dataset, this prior is updated to the “model space posterior” using Bayes’ Rule: $p(\theta|\mathcal{D}) \propto p(\theta)p(\mathcal{D}|\theta)$. This posterior is then used to make predictions.

One method for doing so is to obtain samples from the posterior, $S = \{\theta_1, \dots, \theta_M\} \sim p(\theta|\mathcal{D})$ through a method such as *Markov Chain Monte Carlo (MCMC)*, and use them to approximate the “posterior predictive distribution”:

$$p(y|x, \mathcal{D}) = \int p(y|x, \theta)p(\theta|\mathcal{D})d\theta \approx \frac{1}{M} \sum_j p(y|x, \theta_j).$$

Knowledge Distillation

The goal of knowledge distillation is to compress a given larger “teacher” model $m_{\theta_T}(x)$ into a smaller “student” model $m_{\theta_S}(x)$ that matches its predictions on a shared data distribution (Hinton, Vinyals, and Dean 2015).

A number of FL techniques assume the server has access to a public (unlabelled) dataset \mathcal{U} that serves as the distillation dataset (Guha, Talwalkar, and Smith 2019; Li, He, and Song 2021; Chen and Chao 2021). The student is trained to minimize a loss of the form $\mathcal{L}(\theta_S) = \sum_{x \in \mathcal{U}} l(m_{\theta_S}(x), m_{\theta_T}(x))$. Here, $l(\cdot, \cdot)$ measures the discrepancy between predictions. For classification, this discrepancy can be measured by the Kullback-Leibler divergence between the class distributions, while for regression tasks, it can be measured by the mean-squared error.

Related Work

Bayesian Techniques in FL

Bayesian learning offers FL techniques the advantage of better performance in heterogeneous settings. As argued by Al-Shedivat et al. (2021), FedAvg can be thought of as a technique that obtains the mean of the global posterior, if each local posterior is approximated as a Gaussian with the identity as the covariance matrix. Thus FedAvg implicitly assumes a form of homogeneity, which may not be practical. Bayesian techniques offer the ability to remedy this by using more realistic approximations to the local posteriors.

Existing Bayesian FL techniques focus on approximating the global model space posterior $p(\theta|\mathcal{D})$ from the local model space posteriors $p(\theta|\mathcal{D}_i)$ (Neiswanger, Wang, and Xing 2014; Al-Shedivat et al. 2021).

Embarrassingly Parallel MCMC (EP MCMC, Neiswanger, Wang, and Xing 2014) implements Bayesian inference by drawing MCMC samples from each local posterior (with a corrective factor from the prior), and then estimating the local densities either as Gaussians or with a kernel density estimator. These local densities are then aggregated via multiplication (again with a prior corrective

factor) to obtain an approximation for the global model space posterior. This global density is then sampled to obtain the desired posterior samples. It is worth noting that the original work was not designed for use with neural networks, and the memory costs associated with the method make it intractable for this setting. For instance, when approximating the local posteriors as Gaussians and aggregating them, a computational cost of $O(d^3)$ is required for inverting the covariance matrices, where d is the number of neural network parameters. This method is notable however for operating with only a single communication round.

Federated Posterior Averaging (Al-Shedivat et al. 2021) is similar to EP MCMC, except that it approximates the local posteriors as Gaussians, and devises an efficient iterative algorithm for aggregating the local posteriors (with cost linear in the number of parameters). However, the method requires multiple rounds of communication to be effective.

The main issue with both these techniques is that they require some approximation of the global model space posterior (e.g., in the form of a Gaussian), which can often be inaccurate when the number of model parameters is large. Such approximations are especially poor for neural network models, where the model space posterior is known to be multimodal (Pourzanjani, Jiang, and Petzold 2017).

The **Bayesian Committee Machine** (BCM, Tresp 2000; Deisenroth and Ng 2015) can be thought of as an aggregation method that combines low dimensional predictive posteriors $p(y|x, \mathcal{D}_i)$, rather than the parameter posteriors $p(\theta|\mathcal{D}_i)$. It aggregates as

$$p(y|x, \mathcal{D}) = \frac{1}{p^{m-1}(y|x)} \prod_i p(y|x, \mathcal{D}_i). \quad (1)$$

This formula is correct assuming that the data shards \mathcal{D}_i are independent from each other (and conditionally independent given the single query point (x, y)). We note that this can be satisfied if the data shards form clusters in input space, in other words, if the data is heterogeneous in a certain way.

The BCM results in an ensemble model over the local Bayesian samples. The advantage of this method is that predictive posteriors are much simpler, and lend themselves to, for instance, Gaussian approximations, without sacrificing accuracy. On the other hand, the aggregation formula is no longer exact, but only approximately true, which causes other inaccuracies (specifically, in calibration). The BCM is the starting point for our proposed method.

The **Generalized Robust BCM** (Liu et al. 2018) proposes corrections to the calibration shortcomings of the BCM in regression. Our work analyzes the calibration under different conditions, and extends the analysis to the classification setting. In addition, the proposed correction algorithm requires sharing a subset of the data to all clients, which is incompatible with FL privacy constraints. In this work, we propose a different correction procedure amenable to FL.

One-Shot Federated Learning

Owing to the importance of efficient communication, several methods have been developed to perform FL within a single communication round. **One-Shot FL** (Guha, Talwalkar, and

Smith 2019), as well as **Federated Learning via Knowledge Transfer (FedKT)** (Li, He, and Song 2021) perform one round training by constructing an ensemble from the client models, and compressing it into a single model using knowledge distillation on a public unlabeled dataset. The methods differ in how they construct the teacher ensemble: whereas “One-Shot FL” averages the client predictions (and was tailored for SVM models), “FedKT” aggregates based on majority voting by local models (and only applies to classification tasks). The latter technique uses a scheme called “consistent voting”, where it uses discrepancies between client votes to determine which clients are uncertain about their predictions, and thus can be ignored in the majority vote. Our work is similar to these approaches, but derives aggregation rules for the local client models using a Bayesian perspective, and is applicable to any type of tasks or predictive models. In contrast to our method, these existing techniques do not prioritize well-calibrated predictions.

Analysis of BCM Calibration

As mentioned before, the BCM aggregation requires the independence assumption of the data shards. We analyze when and how this approximation fails: in that it can produce an overconfident global model.

We can analyze the BCM equation in the context of Gaussian process (GP) regression, since it allows us to explicitly calculate the predictive mean and variance.

We assume a smooth, isotropic kernel function $k(x, x') = k(\|x - x'\|)$. We assume a model with Gaussian noise, i.e., $y = f(x) + \epsilon$ where $\epsilon \sim \mathcal{N}(0, \sigma_o^2)$ and σ_o^2 is the “observation variance” in the predictions (in other words, the aleatoric uncertainty in our predictions).

In this setting, the BCM predictive posteriors can be approximated as Gaussians $p(y|x, \mathcal{D}_i) = \mathcal{N}(\mu_i, \Sigma_i)$ (and with prior $p(y|x) = \mathcal{N}(\mu_p, \Sigma_p)$). The aggregation formula (1) now computes a global mean and covariance:

$$\Sigma_g = \left(\sum_i \Sigma_i^{-1} - (n-1)\Sigma_p^{-1} \right)^{-1}, \quad (2)$$

$$\mu_g = \Sigma_g \left(\sum_i \Sigma_i^{-1} \mu_i - (n-1)\Sigma_p^{-1} \mu_p \right). \quad (3)$$

Suppose the input data-points lie in some bounded region R , i.e., $x_i \in R$ for all i . For the GP we outline two observations about the predictive variance $\sigma^2(x_*)$:

Lemma 1 (Choi and Schervish 2007). *Assume $x^* \in R$. Under some mild conditions on the kernel function, and under the assumption of Gaussian or Laplacian observation noise, as the number of data-points increases $\sigma^2(x^*) \rightarrow \sigma_o^2$ (and in addition, the predictive mean converges to the true function value: $\mu(x^*) \rightarrow f(x^*)$).*

Lemma 2. *Assume $x^* \notin R$, and is sufficiently far away from all training points such that $k(x^*, x_i) \approx 0$ for all $x_i \in R$. Then the predictive variance becomes $\sigma^2(x^*) = \sigma_o^2 + k(x^*, x^*)$, which we refer to as the “prior variance” σ_p^2 . Also, the predictive mean becomes $\mu(x^*) = 0$.*

The proof for this lemma, and all following theorems are in the appendix. Taken together, Lemma 1 and 2 mean that near data, the GP attains the correct mean and variance, and away from it, it reverts to a larger prior variance. This reflects the model’s higher uncertainty on unobserved data.

Equipped with these observations, we can analyze the case with partitioned data. We assume two idealized partitions for the overall dataset \mathcal{D} into the m shards:

- **Idealized Homogeneous Partition:** each dataset \mathcal{D}_i is sampled from the same distribution as global dataset \mathcal{D} .
- **Idealized Heterogeneous Partition:** all points $x_i \in \mathcal{D}_i$ and $x_j \in \mathcal{D}_j$ are separated ($k(x_i, x_j) \approx 0$). In other words, the local datasets form clusters in the input space.

Our primary result is that the BCM equations under-estimate the predictive variance in the case of the idealized *homogeneous* partition, i.e., the BCM predictions are overconfident:

Theorem 1 (BCM, homogeneous). *If the data is split among m clients in an idealized homogeneous partition, and $x^* \in R$ is a single test point, then the BCM equations (2) underestimate the predictive variance $\sigma_{\text{BCM}}^2(x^*) < \sigma_o^2$ as the size of the data subsets grows ($|\mathcal{D}_i| \rightarrow \infty$).*

On the other hand, in the same setting, for the mean prediction, if we assume a small prior precision $\sigma_p^{-2} \approx 0$, a quick calculation shows that we still obtain accurate results. Starting with the equation for the BCM predictive mean (2) and applying Lemma 1, we obtain:

$$\begin{aligned} \mu_{\text{BCM}}(x^*) &= \sum_i \sigma_i^{-2} \mu_i / \left(\sum_i \sigma_i^{-2} \right) \\ &\approx \sum_i \sigma_o^{-2} f(x^*) / \left(\sum_i \sigma_o^{-2} \right) = f(x^*). \end{aligned}$$

In other words, the predictive mean still recovers the true mean function at x^* . We can also check that the BCM produces a calibrated predictive variance in the case of idealized heterogeneous data, which lines up with the idea that the BCM is accurate under certain *heterogeneous* partitions. Note that this contrasts with proposition 1 of Liu et al. (2018), which asserts that the BCM is overconfident with heterogeneous partitions. The reason for the difference is that they assume i) that the data stays in the same bounded region as $|\mathcal{D}| \rightarrow \infty$, and ii) that the number of partitions grows to ∞ . In that setting, the disjoint partitions clump together as $|\mathcal{D}| \rightarrow \infty$, unlike our setting.

Theorem 2 (BCM, heterogeneous). *If the data is split among m clients in an idealized heterogeneous partition, and $x^* \in R$ is a test point near \mathcal{D}_k , then the BCM equations (2) correctly estimate the predictive variance $\sigma_{\text{BCM}}^2(x^*) = \sigma_o^2$ as the size of the data subsets grows ($|\mathcal{D}_i| \rightarrow \infty$).*

Analyzing the Predictive Mixture Model

From the above we see that the BCM model produces overconfident estimates for the (idealized) homogeneous setting. An alternate model that produces well-calibrated estimates in this setting is the **predictive mixture model**:

$$\sum_i \frac{|\mathcal{D}_i|}{|\mathcal{D}|} p(y|x, \mathcal{D}_i). \quad (4)$$

For GP regression, when each $p(y|x, \mathcal{D}_i) = \mathcal{N}(\mu_i(x), \sigma_i^2(x))$ is Gaussian, this model produces a Gaussian mixture as its predictive distribution. In general, this will not be Gaussian, but we can approximate it with an overall Gaussian prediction $\mathcal{N}(\mu_{\text{mix}}(x), \sigma_{\text{mix}}^2(x))$, matching the mean and variance of the mixture:

$$\mu_{\text{mix}}(x) = \sum_i \frac{|\mathcal{D}_i|}{|\mathcal{D}|} \mu_i(x), \quad (5)$$

$$\sigma_{\text{mix}}^2(x) = \sum_i \frac{|\mathcal{D}_i|}{|\mathcal{D}|} (\sigma_i^2(x) + \mu_i^2(x)) - \mu_{\text{mix}}^2(x). \quad (6)$$

This is because the moments of a mixture are equal to the mixtures of the moments. Unlike the BCM model, for homogeneous data, the mixture model is well-calibrated:

Theorem 3 (mixture model, homogeneous). *If the data is split among m clients in an idealized homogeneous partition, and $x^* \in R$ is a test point in the bounded region of the data, then the predictive mixture equations (5) correctly estimate the predictive variance $\sigma_{\text{mix}}^2(x^*) = \sigma_o^2$ as the size of the data subsets grows ($|\mathcal{D}_i| \rightarrow \infty$).*

However, the drawback of the mixture model is that for the idealized heterogeneous partition, it overestimates the predictive uncertainty:

Theorem 4 (mixture model, heterogeneous). *If the data is split among m clients in an idealized heterogeneous partition, and $x^* \in R$ is a test point near \mathcal{D}_k , then the predictive mixture equations (5) overestimate the predictive variance $\sigma_{\text{mix}}^2(x^*) > \sigma_o^2$ as the data subset size grows ($|\mathcal{D}_i| \rightarrow \infty$).*

Calibration Analysis for Classification

The preceding analysis may be extended to classification (or more general) models under the following assumptions. Below, we denote the true underlying predictive model (with correct aleatoric uncertainty) as $p_T(y|x)$, and the prior predictive model as $p_P(y|x)$. We assume:

1. For x^* in the vicinity of data subset \mathcal{D}_i , and assuming all data subsets are large enough, the local predictive model converges to the true model, i.e. $p(y|x, \mathcal{D}_i) = p_T(y|x)$;
2. For x^* far away from the data subset \mathcal{D}_i , the local predictor outputs the prior model: $p(y|x, \mathcal{D}_i) = p_P(y|x)$;
3. The prior model $p_P(y|x)$ has higher variance, uncertainty, and/or entropy than the true model $p_T(y|x)$.

All these are desiderata for well-calibrated models. The first desideratum ensures that the model is accurate for in-domain points and the second desideratum ensures that the model reverts to predictions it would have made in the absence of data (the prior predictions) for out-of-domain points. The final desideratum ensures that the prior predictions are sufficiently uncertain.

Although some of these assumptions may not be satisfied by some approximate models, they are the ideal goal for a well-calibrated model. For instance, in classification with neural networks using a softmax layer, item 2 may not be true, and it is the goal of multiple methods to correct this defect (Mukhoti et al. 2021). This is why, for analytically simple models such as GP regression, these assumptions hold.

	Heterogeneous Data	Homogeneous Data
BCM	Calibrated	Overconfident
Mixture	Underconfident	Calibrated

Table 1: Summary of Predictive Variance Analysis.

Under the above assumptions of well-calibrated local models, we can extend the analysis to classification. For simplicity we assume $p_P(y|x)$ is a uniform distribution. We let c denote the correct class for the test point x^* .

Theorem 5 (BCM, homogeneous, classification). *Assuming well-calibrated local models, and an idealized homogeneous partition, if $x^* \in R$ is a test point in the bounded region of the data, then BCM predicts the correct class with probability $p_{\text{BCM}}(c|x^*, \mathcal{D}) > p_T(c|x^*, \mathcal{D})$.*

This follows from the fact the local distributions are correct $p(y|x, \mathcal{D}_i) = p_T(y|x)$ and BCM multiplies them together to produce $p_{\text{BCM}}(y|x, \mathcal{D}) \propto p_T^m(y|x)$, which will be more confident than the desired distribution $p_T(y|x)$.

For heterogeneous data, only one of the factors in the product distribution converges to $p_T(y|x)$, while the others still match the uniform prior $p_P(y|x)$. In this case, the product yields an accurate model.

Theorem 6 (BCM, heterogeneous, classification). *Assuming well-calibrated local models, and an idealized heterogeneous partition, if $x^* \in R$ is a single test point near \mathcal{D}_k , then the BCM predicts the correct class with probability $p_{\text{BCM}}(c|x^*, \mathcal{D}) = p_T(c|x^*, \mathcal{D})$.*

Similarly analyzing the mixture model, we see it is underconfident for heterogeneous data:

Theorem 7 (mixture model, heterogeneous, classification). *Assuming well-calibrated local models, and an idealized heterogeneous partition, if $x^* \in R$ is a test point near \mathcal{D}_k , then the mixture predicts the correct class with probability $p_{\text{mix}}(c|x^*, \mathcal{D}) < p_T(c|x^*, \mathcal{D})$.*

On the other hand, for homogeneous data, each predictive distribution is identical, and equal to the true distribution, so that the mixture returns the correct distribution $p_{\text{mix}}(y|x, \mathcal{D}) = \sum_i (|\mathcal{D}_i|/|\mathcal{D}|) p_T(y|x) = p_T(y|x)$.

Theorem 8 (mixture model, homogeneous, classification). *Assuming well-calibrated local models, and an idealized homogeneous partition, if $x^* \in R$ is a test point in the bounded region of the data, then the mixture predicts the correct class with probability $p_{\text{mix}}(c|x^*, \mathcal{D}) = p_T(c|x^*, \mathcal{D})$.*

We summarize the analysis results in Table 1.

Calibrating the Aggregated Model

The preceding analysis motivates us to combine the predictive mixture model (4), and the BCM (henceforth referred to as the ‘‘product model’’) (1). Namely, to obtain the correct calibration in the predictive model, we should interpolate between the mixture, which is accurate for homogeneous data, and the product, which is accurate for heterogeneous data.

Algorithm 1: Distilled β -PredBayes

Input: Client datasets \mathcal{D}_i , sampler `MCMC_sample`

Output: Model θ^*

```

for each client  $i$  do
   $\{\theta\}_i = \text{MCMC\_sample}(\mathcal{D}_i)$  //step 1
  Communicate  $\{\theta\}_i$  to server //step 2
end for
At Server:
 $\hat{p}(y|x, \mathcal{D}_i) = \frac{1}{|\{\theta\}_i|} \sum_{\theta \in \{\theta\}_i} p(y|x, \theta)$  //step 3
 $\hat{p}_\beta(y|x, \mathcal{D}) = \text{Aggregate}(\hat{p}(y|x, \mathcal{D}_i))$  //step 4, Eq. 7
 $\beta^* = \text{argmin}_\beta \sum -\log \hat{p}_\beta(y|x, \mathcal{D})$  //step 5, tune  $\beta$ 
 $\theta^* = \text{Distill}(\hat{p}_{\beta^*}(y|x, \mathcal{D}))$  //step 6
return  $\theta^*$ 

```

For interpolation parameter β , the model, which we refer to as β -Predictive Bayes is:

$$\log p_\beta(y|x, \mathcal{D}) = \beta \log \left(\frac{1}{p(y|x)^{n-1}} \prod_i p(y|x, \mathcal{D}_i) \right) + (7)$$

$$(1 - \beta) \log \left(\sum_i \frac{|\mathcal{D}_i|}{|\mathcal{D}|} p(y|x, \mathcal{D}_i) \right).$$

The case of $\beta = 0.0$ corresponds to the mixture model, and $\beta = 1.0$ corresponds to the product model.

In the case of regression, assuming Gaussian outputs for each local predictive distribution $p(y|x, \mathcal{D}_i)$, and using the mixture approximations (5), the aggregated predictive distribution is approximately a Gaussian $p_\beta(y|x, \mathcal{D}) = \mathcal{N}(\mu_\beta(x), \sigma_\beta^2(x))$ with:

$$\sigma_\beta^{-2}(x) = \beta \cdot \sigma_{\text{prod}}^{-2} + (1 - \beta) \cdot \sigma_{\text{mix}}^{-2}(x) \quad (8)$$

$$\mu_\beta(x) = \sigma_\beta^2(x) (\beta \cdot \sigma_{\text{prod}}^{-2}(x) \mu_{\text{prod}}(x) + (1 - \beta) \cdot \sigma_{\text{mix}}^{-2}(x) \mu_{\text{mix}}(x)).$$

In other words, the inverse-variance (or precision) interpolates between that of the product and mixture distributions.

We learn β by minimizing the negative log-likelihood of a dataset on the server, \mathcal{U} using a gradient-based optimizer.

$$\beta^* = \text{argmin}_\beta \sum_{(x,y) \in \mathcal{U}} -\log p_\beta(y|x, \mathcal{D}). \quad (9)$$

For regression, the Gaussian negative log-likelihood can be used for training (by approximating $p_\beta(y|x, \mathcal{D})$ with (8)).

By tuning β this way, we don’t need to know where we are along the homogeneous-heterogeneous partition spectrum. Since we are only tuning a single scalar parameter, a small dataset may suffice. Furthermore, we can use \mathcal{U} as a distillation dataset to compress the ensemble into one model. The steps for β -predBayes are presented in Algorithm 1.

Experiments

We evaluate β -PredBayes on a number of regression and classification datasets. All tests used 5 clients, with a distillation set composed of 20% of the original training set. Further experimental details, such as the models, and hyperparameters, are included in the appendix.

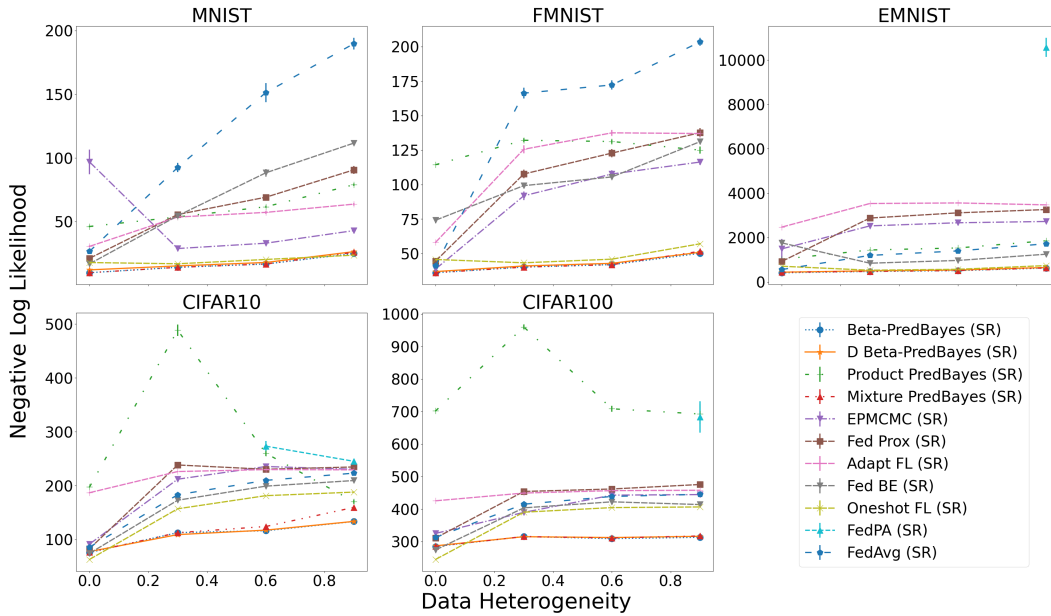


Figure 1: NLL on the classification datasets with increasing heterogeneity (tested with $h \in \{0.0, 0.3, 0.6, 0.9\}$). Averages and standard error over 10 seeds are reported. Omitted values (e.g., for FedPA on EMNIST) denote results where NLL diverged.

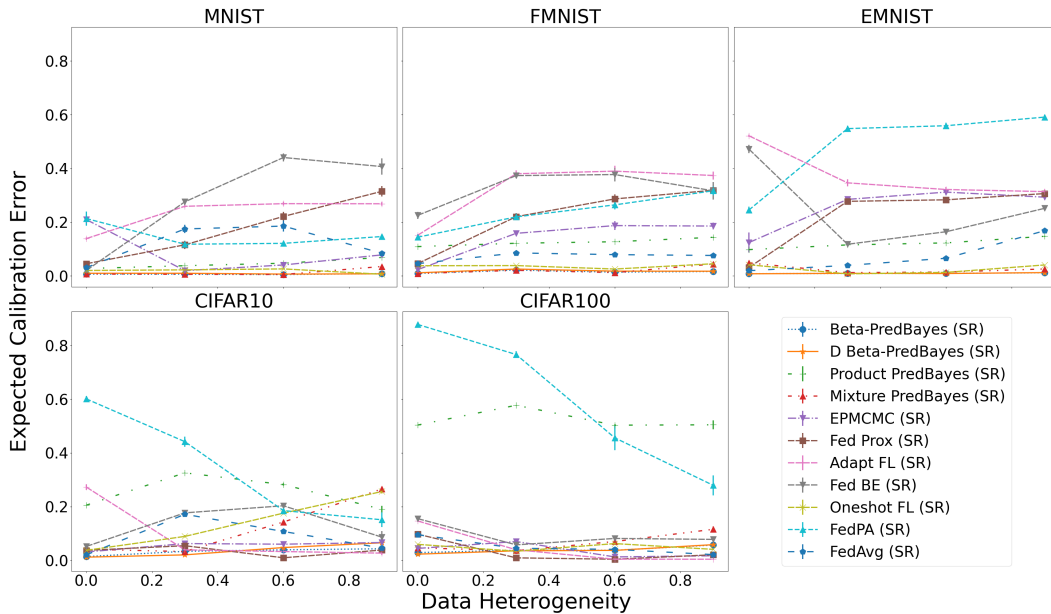


Figure 2: ECE on the classification datasets with increasing heterogeneity (tested with $h \in \{0.0, 0.3, 0.6, 0.9\}$). Averages and standard error over 10 seeds are reported.

Since the goal is to evaluate the calibration of the predictions (not just their accuracy), we measure negative-log-likelihood (NLL) on the test set. This evaluates calibration since NLL should be minimized for the correct prediction probabilities (Guo et al. 2017). For classification, we may also use the expected calibration error (ECE), which measures the difference between the confidence (probabilities) assigned to predicted classes, and the fraction that are pre-

dicted correctly (see the appendix for more details).

Classification Dataset Setup

The method was evaluated for classification on the following datasets: MNIST (Lecun et al. 1998), Fashion MNIST (Xiao, Rasul, and Vollgraf 2017), EMNIST (Cohen et al. 2017) (using a split with 62 classes), CIFAR10 and CIFAR100 (Krizhevsky, Hinton et al. 2009). Each of these

Method	Air Quality	Bike	Wine Quality	Real Estate	Forest Fire
EP-MCMC	11.20 ± 0.42	1.92 ± 0.12	2.93 ± 0.09	1.90 ± 0.29	1.90 ± 0.09
FedBE	6.34 ± 0.03	0.94 ± 0.05	2.13 ± 0.01	0.49 ± 0.01	1.39 ± 0.01
OneshotFL	6.69 ± 0.05	0.98 ± 0.05	2.17 ± 0.01	0.53 ± 0.01	1.39 ± 0.01
Mixture	9.45 ± 0.02	1.22 ± 0.02	2.57 ± 0.03	0.65 ± 0.01	1.39 ± 0.01
Product	10.17 ± 0.08	1.47 ± 0.02	3.06 ± 0.10	0.83 ± 0.03	2.57 ± 0.03
β -PredBayes (ours)	7.05 ± 0.04	0.95 ± 0.01	2.05 ± 0.05	0.53 ± 0.03	1.55 ± 0.01*
D- β -PredBayes (ours)	4.53 ± 0.12	0.32 ± 0.05	1.34 ± 0.03	0.18 ± 0.02	1.56 ± 0.03

Table 2: Average NLL (\pm standard error) on regression datasets, based on 10 seeds. Lower is better. All values are statistically significant relative to D- β -PredBayes ($p < 5\%$) according to the Wilcoxon signed-rank test, except those marked with *.

datasets is distributed among clients based on a “heterogeneity parameter”, h . For $h = 0$, the data is split uniformly among the clients, and is thus homogeneously distributed. For $h = 1$, the data is sorted by class before being split among clients. This means that each client observes data from different classes with little overlap. In this sense, $h = 1$ represents the “fully heterogeneous” setting. For $0 < h < 1$, data from the above extremes is mixed: a fraction of size h of the homogeneous data shard is replaced with the fully heterogeneous data for each client.

Classification Results

We evaluate β -PredBayes, in both its distilled (listed as D- β -PredBayes) and non-distilled forms, along with the product (BCM) and mixture models. We also evaluate several baselines, run for one communication round. The results for NLL over different settings of the heterogeneity are plotted in Figure 1, while ECE results are shown in Figure 2.

For the NLL, the β -PredBayes model (and its distilled variant) perform best (least NLL), followed by the Mixture model, on the tested datasets. As heterogeneity increases, the NLL loss generally increases for other methods, while it largely remains stable for β -PredBayes. For the ECE, as heterogeneity increases, the metric jumps more erratically for some methods (which may be due ECE’s sensitivity to hyperparameters like bin count). But the overall trend is still that β -PredBayes performs best (with lowest ECE), followed by its distilled variant, followed by the mixture model. We note that for ECE, β -PredBayes outperforms the mixture model for high heterogeneity setting ($h = 0.9$), which is expected from the analysis that predicted that mixture models would not be well-calibrated in this setting.

Regression Dataset Setup

The regression datasets used for evaluation include: the “wine quality” (Cortez et al. 2009), “air quality” (De Vito et al. 2008), “forest fire” (Cortez and Morais 2007), “real estate” (Yeh and Hsu 2018), and “bike rental” (Fanaee-T and Gama 2013) datasets from the UCI repository (Dua and Graff 2017). These datasets were sorted according to certain features (such as the date, for “airquality”), then split among clients, to simulate heterogeneous data.

Regression Results

For regression, β -predBayes was evaluated, along with the product and mixture models as well as other baselines that use an ensemble for predictive inference (since these yield a predictive variance). The resulting (Gaussian) NLL for heterogeneous data are shown in Table 2. For all datasets except “Forest Fire,” the distilled form of β -PredBayes performs best in NLL. Note that it is possible that the distilled version outperforms the non-distilled one due to some regularization effect of having a smaller model.

Conclusion and Future Work

This work presented β -Predictive Bayes, an algorithm that aggregates local Bayesian predictive posteriors using a tunable parameter β , and then distills the resulting model. The parameter β enables accurate calibration performance. Owing to its Bayesian nature, the method is also well suited for heterogeneous FL, and requires only a single communication round. We perform empirical evaluation on various classification and regression datasets to show the competence of our methods in heterogeneous settings. Our work reinforces that Bayesian learning can provide well-calibrated models in heterogeneous settings with efficient communication.

Some future directions to improve our work include:

- **Personalized FL.** Our proposed approach learns a single global model for all clients. When clients have different class-conditional distributions $p(y|x)$, a personalized model can be more desirable for each client.
- **Privacy.** The proposed method transmits weight samples (obtained by MCMC) to the server. Although no data is shared, information about the data could be leaked via the model weights unless a theoretically proven private mechanism is used or encryption is applied. This trait is shared by other techniques (e.g., FedAvg, FedPA). It would be interesting to explore rigorously private methods, such as differentially private sampling mechanisms (Dimitrakakis et al. 2017).
- **Server dataset.** β -Predictive Bayes uses public data stored at the server for distillation and tuning β . When a public dataset is not available, it may need to be synthetically generated. It would be better to avoid this generation and to develop a data-free technique for distillation and learning β in the future.

Acknowledgments

Resources used in preparing this research at the University of Waterloo were provided by Huawei Canada, the Natural Sciences and Engineering Research Council of Canada, the Province of Ontario, the Government of Canada through CI-FAR, and companies sponsoring the Vector Institute.

References

- Al-Shedivat, M.; Gillenwater, J.; Xing, E.; and Rostamizadeh, A. 2021. Federated Learning via Posterior Averaging: A New Perspective and Practical Algorithms. In *International Conference on Learning Representations*.
- Bishop, C. M. 2006. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Berlin, Heidelberg: Springer-Verlag. ISBN 0387310738.
- Chen, H.-Y.; and Chao, W.-L. 2021. Fed{BE}: Making Bayesian Model Ensemble Applicable to Federated Learning. In *International Conference on Learning Representations*.
- Choi, T.; and Schervish, M. J. 2007. On posterior consistency in nonparametric regression problems. *Journal of Multivariate Analysis*, 98(10): 1969–1987.
- Cohen, G.; Afshar, S.; Tapson, J.; and Van Schaik, A. 2017. EMNIST: Extending MNIST to handwritten letters. In *2017 international joint conference on neural networks (IJCNN)*, 2921–2926. IEEE.
- Cortez, P.; Cerdeira, A.; Almeida, F.; Matos, T.; and Reis, J. 2009. Modeling wine preferences by data mining from physicochemical properties. *Decision support systems*, 47(4): 547–553.
- Cortez, P.; and Morais, A. d. J. R. 2007. A data mining approach to predict forest fires using meteorological data.
- De Vito, S.; Massera, E.; Piga, M.; Martinotto, L.; and Di Francia, G. 2008. On field calibration of an electronic nose for benzene estimation in an urban pollution monitoring scenario. *Sensors and Actuators B: Chemical*, 129(2): 750–757.
- Deisenroth, M. P.; and Ng, J. W. 2015. Robust Bayesian Committee Machine for Large-Scale Gaussian Processes. In *Large-Scale Kernel Machines Workshop at ICML 2015*.
- Dimitrakakis, C.; Nelson, B.; Zhang, Z.; Mitrokovtsa, A.; and Rubinstein, B. I. 2017. Differential privacy for Bayesian inference through posterior sampling. *Journal of machine learning research*, 18(11): 1–39.
- Dua, D.; and Graff, C. 2017. UCI Machine Learning Repository.
- Fanaee-T, H.; and Gama, J. 2013. Event labeling combining ensemble detectors and background knowledge. *Progress in Artificial Intelligence*, 1–15.
- Guha, N.; Talwalkar, A.; and Smith, V. 2019. One-Shot Federated Learning.
- Guo, C.; Pleiss, G.; Sun, Y.; and Weinberger, K. Q. 2017. On calibration of modern neural networks.
- Hinton, G. E.; Vinyals, O.; and Dean, J. 2015. Distilling the Knowledge in a Neural Network. *ArXiv*, abs/1503.02531.
- Kingma, D. P.; and Ba, J. 2015. Adam: A Method for Stochastic Optimization. In Bengio, Y.; and LeCun, Y., eds., *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Krizhevsky, A.; Hinton, G.; et al. 2009. Learning multiple layers of features from tiny images.
- Lecun, Y.; Bottou, L.; Bengio, Y.; and Haffner, P. 1998. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11): 2278–2324.
- Li, Q.; He, B.; and Song, D. 2021. Practical One-Shot Federated Learning for Cross-Silo Setting. In Zhou, Z.-H., ed., *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, 1484–1490. International Joint Conferences on Artificial Intelligence Organization. Main Track.
- Li, T.; Sahu, A. K.; Zaheer, M.; Sanjabi, M.; Talwalkar, A.; and Smith, V. 2020a. Federated optimization in heterogeneous networks. *Proceedings of Machine Learning and Systems*, 2: 429–450.
- Li, X.; JIANG, M.; Zhang, X.; Kamp, M.; and Dou, Q. 2020b. FedBN: Federated Learning on Non-IID Features via Local Batch Normalization. In *International Conference on Learning Representations*.
- Liu, H.; Cai, J.; Wang, Y.; and Ong, Y. S. 2018. Generalized Robust Bayesian Committee Machine for Large-scale Gaussian Process Regression. In *ICML*, 3131–3140.
- McMahan, B.; Moore, E.; Ramage, D.; Hampson, S.; and Arcas, B. A. y. 2017. Communication-Efficient Learning of Deep Networks from Decentralized Data. In Singh, A.; and Zhu, J., eds., *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, 1273–1282. PMLR.
- Mohri, M.; Sivek, G.; and Suresh, A. T. 2019. Agnostic federated learning. In *International Conference on Machine Learning*, 4615–4625. PMLR.
- Mukhoti, J.; Kirsch, A.; van Amersfoort, J.; Torr, P. H.; and Gal, Y. 2021. Deterministic Neural Networks with Appropriate Inductive Biases Capture Epistemic and Aleatoric Uncertainty. *arXiv preprint arXiv:2102.11582*.
- Neiswanger, W.; Wang, C.; and Xing, E. P. 2014. Asymptotically Exact, Embarrassingly Parallel MCMC. In *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence*, UAI'14, 623–632. Arlington, Virginia, USA: AUAI Press. ISBN 9780974903910.
- Pourzanjani, A. A.; Jiang, R. M.; and Petzold, L. R. 2017. Improving the identifiability of neural networks for Bayesian inference. In *NeurIPS Workshop on Bayesian Deep Learning*.
- Rasmussen, C. E.; and Williams, C. K. I. 2006. *Gaussian processes for machine learning*. Adaptive computation and machine learning. MIT Press. ISBN 026218253X.
- Reddi, S. J.; Charles, Z.; Zaheer, M.; Garrett, Z.; Rush, K.; Konečný, J.; Kumar, S.; and McMahan, H. B. 2021. Adaptive Federated Optimization. In *International Conference on Learning Representations*.

Tresp, V. 2000. A Bayesian committee machine. *Neural computation*, 12(11): 2719–2741.

Wang, H.; Yurochkin, M.; Sun, Y.; Papailiopoulos, D.; and Khazaeni, Y. 2019. Federated Learning with Matched Averaging. In *International Conference on Learning Representations*.

Wang, J.; Liu, Q.; Liang, H.; Joshi, G.; and Poor, H. V. 2020. Tackling the objective inconsistency problem in heterogeneous federated optimization. *Advances in neural information processing systems*, 33: 7611–7623.

Woodworth, B. E.; Patel, K. K.; and Srebro, N. 2020. Mini-batch vs Local SGD for Heterogeneous Distributed Learning. *CoRR*, abs/2006.04735.

Xiao, H.; Rasul, K.; and Vollgraf, R. 2017. Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms.

Yeh, I.-C.; and Hsu, T.-K. 2018. Building real estate valuation models with comparative approach through case-based reasoning. *Applied Soft Computing*, 65: 260–271.

Zhang, R.; Li, C.; Zhang, J.; Chen, C.; and Wilson, A. G. 2020. Cyclical Stochastic Gradient MCMC for Bayesian Deep Learning. *International Conference on Learning Representations*.

Zhao, Y.; Li, M.; Lai, L.; Suda, N.; Civin, D.; and Chandra, V. 2018. Federated learning with non-iid data. *arXiv preprint arXiv:1806.00582*.

Zheng, Z.; Zhou, Y.; Sun, Y.; Wang, Z.; Liu, B.; and Li, K. 2022. Applications of federated learning in smart cities: recent advances, taxonomy, and open challenges. *Connection Science*, 34(1): 1–28.

A Proofs of Theorems

We provide proofs for the theorems in the main text. Recall the BCM aggregation equation as:

$$p(y|x, \mathcal{D}) = \frac{1}{p^{m-1}(y|x)} \prod_i p(y|x, \mathcal{D}_i)$$

which, in the case of Gaussian local predictive posteriors $p(y|x, \mathcal{D}_i) = \mathcal{N}(\mu_i, \Sigma_i)$ simplifies to computing the Gaussian global predictions $p(y|x, \mathcal{D}) = \mathcal{N}(\mu_g, \Sigma_g)$:

$$\Sigma_g = \left(\sum_i \Sigma_i^{-1} - (n-1)\Sigma_p^{-1} \right)^{-1}$$

$$\mu_g = \Sigma_g \left(\sum_i \Sigma_i^{-1} \mu_i - (n-1)\Sigma_p^{-1} \mu_p \right).$$

The mixture aggregation is:

$$\sum_i \frac{|\mathcal{D}_i|}{|\mathcal{D}|} p(y|x, \mathcal{D}_i)$$

which we approximate with an overall Gaussian prediction $\mathcal{N}(\mu_{\text{mix}}(x), \sigma_{\text{mix}}^2(x))$:

$$\mu_{\text{mix}}(x) = \sum_i \frac{|\mathcal{D}_i|}{|\mathcal{D}|} \mu_i(x)$$

$$\sigma_{\text{mix}}^2(x) = \sum_i \frac{|\mathcal{D}_i|}{|\mathcal{D}|} (\sigma_i^2(x) + \mu_i^2(x)) - \mu_{\text{mix}}^2(x).$$

Gaussian Process Analysis

We assume a smooth, isotropic kernel function $k(x, x') = k(\|x - x'\|)$. We assume a model with Gaussian noise, i.e., $y = f(x) + \epsilon$ where $\epsilon \sim \mathcal{N}(0, \sigma_o^2)$.

On some dataset $\mathcal{D} = (\mathbf{X}, \mathbf{y})$, and with test point x^* , the Gaussian process inference equations predict a mean and variance for $p(y^*|x^*, \mathcal{D}_i) = \mathcal{N}(\mu(x^*), \sigma^2(x^*))$. If we denote $\mathbf{k}_* = [k(x^*, x_1), \dots, k(x^*, x_{|\mathcal{D}|})]^\top$ the vector of kernel evaluations between the test point x^* and the points in the dataset \mathcal{D} , and \mathbf{K} as the kernel matrix (where $K_{i,j} = k(x_i, x_j)$ for $x_i, x_j \in \mathcal{D}$) then the inference equations are (Rasmussen and Williams 2006):

$$\mu(x^*) = \mathbf{k}_*^\top (\mathbf{K} + \sigma_o^2 I)^{-1} \mathbf{y} \quad (10)$$

$$\sigma^2(x^*) = \sigma_o^2 + k(x^*, x^*) - \mathbf{k}_*^\top (\mathbf{K} + \sigma_o^2 I)^{-1} \mathbf{k}_* \quad (11)$$

We assume the data lies in some bounded set R . Starting with the lemmas in the main text:

Lemma 1 (Choi and Schervish 2007). *Assume $x^* \in R$. Under some mild conditions on the kernel function, and under the assumption of Gaussian or Laplacian observation noise, as the number of data-points increases $\sigma^2(x^*) \rightarrow \sigma_o^2$ (and in addition, the predictive mean converges to the true function value: $\mu(x^*) \rightarrow f(x^*)$).*

Lemma 2. *Assume $x^* \notin R$, and is sufficiently far away from all training points such that $k(x^*, x_i) \approx 0$ for all $x_i \in R$. Then the predictive variance becomes $\sigma^2(x^*) = \sigma_o^2 + k(x^*, x^*)$, which we refer to as the “prior variance” σ_p^2 . Also, the predictive mean becomes $\mu(x^*) = 0$.*

Proof. The inference equation 10 for the predictive variance reads:

$$\sigma^2(x^*) \approx \sigma_o^2 + k(x^*, x^*) - \mathbf{0}^\top (\mathbf{K} + \sigma_o^2 I)^{-1} \mathbf{0}$$

$$\approx \sigma_o^2 + k(x^*, x^*)$$

Similarly, the predictive mean is $\mu(x^*) \approx \mathbf{0}^\top (\mathbf{K} + \sigma_o^2 I)^{-1} \mathbf{y} = 0$. \square

BCM The primary theorems establishing the calibration performance of the BCM and mixture models are:

Theorem 1 (BCM, homogeneous). *If the data is split among m clients in an idealized homogeneous partition, and $x^* \in R$ is a single test point, then the BCM equations (2) underestimate the predictive variance $\sigma_{\text{BCM}}^2(x^*) < \sigma_o^2$ as the size of the data subsets grows ($|\mathcal{D}_i| \rightarrow \infty$).*

Proof. According to the BCM aggregation equation for regression:

$$\sigma_{\text{BCM}}^{-2}(x^*) = \sum_i \sigma_i^{-2}(x^*) - (m-1)\sigma_p^{-2}(x^*)$$

Where $\sigma_i^{-2}(x^*)$ is the inverse of the predictive variance output by the GP trained on \mathcal{D}_i . As the size of the dataset \mathcal{D}_i increases, from Lemma 1 we know that $\sigma_i(x^*)$ will converge to σ_o^2 . On the other hand $\sigma_p^2 = k(x^*, x^*) + \sigma_o^2$. Combining these facts (in the limit of increasing data points):

$$\begin{aligned} \sigma_{\text{BCM}}^{-2}(x^*) &= m\sigma_o^{-2} - (m-1)(\sigma_o^2 + k(x^*, x^*))^{-1} \\ &> m\sigma_o^{-2} - (m-1)\sigma_o^{-2} \\ &= \sigma_o^{-2} \end{aligned}$$

Where the second inequality follows from the fact that $k(x^*, x^*) > 0$ (for a positive kernel function).

This implies $\sigma_{\text{BCM}}^2 < \sigma_o^2$ as claimed.

Furthermore, for a kernel function with a large prior precision $\sigma_p^{-2} \approx 0$ (which is typical in practice when using an uninformative prior), we see that:

$$\begin{aligned} \sigma_{\text{BCM}}^{-2} &= m\sigma_o^{-2} - (m-1)\sigma_p^{-2} \\ &\approx m\sigma_o^{-2}. \end{aligned}$$

So we end up underestimating the predictive variance by a factor of m , for m clients. \square

Theorem 2 (BCM, heterogeneous). *If the data is split among m clients in an idealized heterogeneous partition, and $x^* \in R$ is a test point near \mathcal{D}_k , then the BCM equations (2) correctly estimate the predictive variance $\sigma_{\text{BCM}}^2(x^*) = \sigma_o^2$ as the size of the data subsets grows ($|\mathcal{D}_i| \rightarrow \infty$).*

Proof. The BCM aggregation equation for regression is:

$$\sigma_{\text{BCM}}^{-2}(x^*) = \sum_i \sigma_i^{-2}(x^*) - (m-1)\sigma_p^{-2}(x^*)$$

Since x^* is in the vicinity of some data subset \mathcal{D}_k , we can apply Lemma 1 to get $\sigma_k^{-2}(x^*) = \sigma_o^{-2}$. On the other hand, since the data subsets are split up heterogeneously, the test point x^* is far from the other data sub-sets such that we can apply Lemma 2 to obtain $\sigma_i^{-2}(x^*) = \sigma_p^{-2}(x^*)$ for all $i \neq k$. Plugging these into the above expression yields:

$$\begin{aligned} \sigma_{\text{BCM}}^{-2}(x^*) &= \sigma_k^{-2}(x^*) + \sum_{i \neq k} \sigma_i^{-2}(x^*) - (m-1)\sigma_p^{-2}(x^*) \\ &= \sigma_o^{-2} + (m-1)\sigma_p^{-2}(x^*) - (m-1)\sigma_p^{-2}(x^*) \\ &= \sigma_o^{-2}. \end{aligned}$$

\square

Mixture Model The following theorems analyze the calibration performance of the mixture model.

Theorem 3 (mixture model, homogeneous). *If the data is split among m clients in an idealized homogeneous partition, and $x^* \in R$ is a test point in the bounded region of the data, then the predictive mixture equations (5) correctly estimate the predictive variance $\sigma_{\text{mix}}^2(x^*) = \sigma_o^2$ as the size of the data subsets grows ($|\mathcal{D}_i| \rightarrow \infty$).*

Proof. In this setting, appealing to Lemma 1, each local predictor outputs $\mathcal{N}(\mu_i(x^*) = f(x^*), \sigma_i^2(x^*) = \sigma_o^2)$. Therefore we have:

$$\begin{aligned} \sigma_{\text{mix}}^2 &= \sum_i \frac{|\mathcal{D}_i|}{|\mathcal{D}|} (\sigma_i^2(x^*) + \mu_i^2(x^*)) - \mu_{\text{mix}}^2(x^*) \\ &= \sum_i \frac{|\mathcal{D}_i|}{|\mathcal{D}|} (\sigma_o^2 + f^2(x^*)) - \left(\sum_i \frac{|\mathcal{D}_i|}{|\mathcal{D}|} f(x^*) \right)^2 \\ &= \sigma_o^2 + f^2(x^*) - f^2(x^*) \\ &= \sigma_o^2. \end{aligned}$$

\square

Theorem 4 (mixture model, heterogeneous). *If the data is split among m clients in an idealized heterogeneous partition, and $x^* \in R$ is a test point near \mathcal{D}_k , then the predictive mixture equations (5) overestimate the predictive variance $\sigma_{\text{mix}}^2(x^*) > \sigma_o^2$ as the data subset size grows ($|\mathcal{D}_i| \rightarrow \infty$).*

Proof. The mixture models predictive variance is:

$$\sigma_{\text{mix}}^2(x^*) = \sum_i \frac{|\mathcal{D}_i|}{|\mathcal{D}|} (\sigma_i^2(x^*) + \mu_i^2(x^*)) - \mu_{\text{mix}}^2(x^*)$$

Since x^* lies in the vicinity of \mathcal{D}_k , we have (applying Lemma 1) that $\sigma_k^2(x^*) = \sigma_o^2$ (and $\mu_k(x^*) = f(x^*)$), while from Lemma 2 we know $\sigma_i^2(x^*) = \sigma_p^2(x^*) > \sigma_o^2$ (and $\mu_i(x^*) = 0$) for all $i \neq k$.

$$\begin{aligned} \sigma_{\text{mix}}^2(x^*) &= \frac{|\mathcal{D}_k|}{|\mathcal{D}|} (\sigma_o^2 + f^2(x^*)) \\ &\quad + \sum_{i \neq k} \frac{|\mathcal{D}_i|}{|\mathcal{D}|} (\sigma_p^2(x^*)) - \mu_{\text{mix}}^2(x^*) \\ &= \frac{|\mathcal{D}_k|}{|\mathcal{D}|} \sigma_o^2 + \sigma_p^2(x^*) \left(1 - \frac{|\mathcal{D}_k|}{|\mathcal{D}|}\right) + \frac{|\mathcal{D}_k|}{|\mathcal{D}|} f^2(x^*) \\ &\quad - \frac{|\mathcal{D}_k|^2}{|\mathcal{D}|^2} f^2(x^*) \\ &= \frac{|\mathcal{D}_k|}{|\mathcal{D}|} \sigma_o^2 + \sigma_p^2(x^*) \left(1 - \frac{|\mathcal{D}_k|}{|\mathcal{D}|}\right) \\ &\quad + \left(\frac{|\mathcal{D}_k|}{|\mathcal{D}|} - \frac{|\mathcal{D}_k|^2}{|\mathcal{D}|^2} \right) f^2(x^*) \\ &> \frac{|\mathcal{D}_k|}{|\mathcal{D}|} \sigma_o^2 + (\sigma_o^2 + k(x^*, x^*)) \left(1 - \frac{|\mathcal{D}_k|}{|\mathcal{D}|}\right) \\ &= \sigma_o^2 + k(x^*, x^*) \left(1 - \frac{|\mathcal{D}_k|}{|\mathcal{D}|}\right) \\ &> \sigma_o^2. \end{aligned}$$

In other words, we overestimate σ_o^2 by a constant of at least $k(x^*, x^*)\left(1 - \frac{|\mathcal{D}_k|}{|\mathcal{D}|}\right)$. \square

We note that in this case, the predictive mean of the model is also inaccurate since $\mu_{\text{mix}}(x^*) = \frac{|\mathcal{D}_k|}{|\mathcal{D}|}f(x^*) \neq f(x^*)$.

Classification Analysis

For classification, the following theorems determine the calibration performance of the BCM and mixture models (assuming a uniform prior). We make use of the fact that for the predicted class $p_T(c|x^*) > p_T(i|x^*)$, for $i \neq c$ (ie. it has the highest probability).

Theorem 5 (BCM, homogeneous, classification). *Assuming well-calibrated local models, and an idealized homogeneous partition, if $x^* \in R$ is a test point in the bounded region of the data, then BCM predicts the correct class with probability $p_{\text{BCM}}(c|x^*, \mathcal{D}) > p_T(c|x^*, \mathcal{D})$.*

Proof. Let $p_i = p_T(y = i|x^*, \mathcal{D})$ denote the true class probabilities. Since the local models are well-calibrated, each local model predicts $p(c|x^*, \mathcal{D}_i) = p_c$. The BCM prediction for the correct class is $p_{\text{BCM}}(c|x^*, \mathcal{D}) \propto p_c^m$. We have:

$$\begin{aligned} p_{\text{BCM}}(c|x^*, \mathcal{D}) &= p_c^m / \left(\sum_i p_i^m\right) \\ &= 1 / \left(\sum_i \left(\frac{p_i}{p_c}\right)^m\right) \\ &> 1 / \left(\sum_i \frac{p_i}{p_c}\right) \\ &= p_c \\ &= p_T(c|x^*, \mathcal{D}). \end{aligned}$$

The inequality in line 3 above follows from the fact that $p_i < p_c$ for the correct class c , so that $p_i/p_c < 1$ and $(p_i/p_c)^m < p_i/p_c$. \square

Theorem 6 (BCM, heterogeneous, classification). *Assuming well-calibrated local models, and an idealized heterogeneous partition, if $x^* \in R$ is a single test point near \mathcal{D}_k , then the BCM predicts the correct class with probability $p_{\text{BCM}}(c|x^*, \mathcal{D}) = p_T(c|x^*, \mathcal{D})$.*

Proof. We have that for x^* near \mathcal{D}_k , $p(c|x^*, \mathcal{D}_k) = p_T(c|x^*, \mathcal{D})$ while for other shards $i \neq k$ $p(c|x^*, \mathcal{D}_i) = p_P(c|x^*) = 1/K$:

$$\begin{aligned} p_{\text{BCM}}(c|x^*, \mathcal{D}) &= \frac{1}{p^{m-1}(c|x^*)} \prod_i p(c|x^*, \mathcal{D}_i) \\ &= K^{m-1} p_T(c|x^*, \mathcal{D}) \prod_{i \neq k} \frac{1}{K} \\ &= \frac{K^{m-1}}{K^{m-1}} p_T(c|x^*, \mathcal{D}) \\ &= p_T(c|x^*, \mathcal{D}). \end{aligned}$$

\square

Theorem 7 (mixture model, heterogeneous, classification). *Assuming well-calibrated local models, and an idealized heterogeneous partition, if $x^* \in R$ is a test point near \mathcal{D}_k , then the mixture predicts the correct class with probability $p_{\text{mix}}(c|x^*, \mathcal{D}) < p_T(c|x^*, \mathcal{D})$.*

Proof. We have that for x^* near \mathcal{D}_k , $p(c|x^*, \mathcal{D}_k) = p_T(c|x^*, \mathcal{D})$ while for other shards $i \neq k$ $p(c|x^*, \mathcal{D}_i) = p_P(c|x^*) = 1/K$. Letting $p_T(c|x^*, \mathcal{D}) = p_c$ and noting that $p_c > \frac{1}{K}$ (since otherwise the data point x^* has no correct label):

$$\begin{aligned} p_{\text{mix}}(c|x^*, \mathcal{D}) &= \sum_i \frac{|\mathcal{D}_i|}{|\mathcal{D}|} p(c|x^*, \mathcal{D}_i) \\ &= \frac{|\mathcal{D}_k|}{|\mathcal{D}|} p_T(c|x^*, \mathcal{D}) + \left(1 - \frac{|\mathcal{D}_k|}{|\mathcal{D}|}\right) p_P(c|x^*) \\ &= \frac{|\mathcal{D}_k|}{|\mathcal{D}|} p_c + \left(1 - \frac{|\mathcal{D}_k|}{|\mathcal{D}|}\right) \frac{1}{K} \\ &< \frac{|\mathcal{D}_k|}{|\mathcal{D}|} p_c + \left(1 - \frac{|\mathcal{D}_k|}{|\mathcal{D}|}\right) p_c \\ &= p_c. \end{aligned}$$

\square

Theorem 8 (mixture model, homogeneous, classification). *Assuming well-calibrated local models, and an idealized homogeneous partition, if $x^* \in R$ is a test point in the bounded region of the data, then the mixture predicts the correct class with probability $p_{\text{mix}}(c|x^*, \mathcal{D}) = p_T(c|x^*, \mathcal{D})$.*

Proof. For each local dataset, $p(c|x^*, \mathcal{D}_i) = p_T(c|x^*, \mathcal{D})$, so that:

$$\begin{aligned} p_{\text{mix}}(c|x^*, \mathcal{D}) &= \sum_i \frac{|\mathcal{D}_i|}{|\mathcal{D}|} p(c|x^*, \mathcal{D}_i) \\ &= \sum_i \frac{|\mathcal{D}_i|}{|\mathcal{D}|} p_T(c|x^*, \mathcal{D}) \\ &= p_T(c|x^*, \mathcal{D}). \end{aligned}$$

\square

B Additional Experimental Details

Additional details for the experiments are provided below.

Hardware, Software, and Randomization Details

The code for experiments was written in the Python language (version 3.8.10), primarily using the Pytorch (version 1.9.0), Numpy (version 1.19.5) and Scipy (version 1.6.2) libraries. Randomization was done by setting seeds for Pytorch and Numpy.

Experiments were carried out on a compute cluster using a single Nvidia GPU (either the T4, or P100).

Expected Calibration Error

Assuming the predicted class probabilities are organized into M equally spaced bins B_1, \dots, B_M on the interval $[0, 1]$, the expected calibration error is calculated as:

$$\sum_{m=1}^M \frac{|B_m|}{|\mathcal{D}|} |\text{acc}(B_m) - \text{conf}(B_m)|, \quad (12)$$

where $\text{acc}(B_m)$ is the average accuracy of the predictions in bin m (the fraction of points in B_m classified correctly) and $\text{conf}(B_m)$ is the average confidence of the predictions in bin m (the average value of $p(c|x^*, \mathcal{D})$, where c denotes the predicted class). For our experiments, we calculated ECE using 15 equally spaced bins.

Models

A two-layer fully connected network with 100 hidden units was used for MNIST, Fashion MNIST, EMNIST, as well as the regression datasets. For the CIFAR10 and CIFAR100 datasets a Convolutional Neural network was used with 3 convolution layers, each followed by “Max Pooling” layers, with a single fully connected layer at the end. For all networks, the ReLU activation function was used between layers. For classification, the output was the predictive distribution over classes, while for regression the output was the mean of the predictive distribution. A Gaussian prior is assumed over the network parameters, $p(\theta) = \mathcal{N}(0, \sigma^2 I)$, with $\sigma = 5e4$.

Baselines

The Federated techniques compared include: Federated Averaging (FedAvg, McMahan et al. 2017, Federated Posterior Averaging (FedPA, Al-Shedivat et al. 2021, Embarrassingly Parallel MCMC (EP MCMC, Neiswanger, Wang, and Xing 2014, FedProx (Li et al. 2020a), Adaptive FL (Reddi et al. 2021), Federated Bayesian Ensemble (FedBE) (Chen and Chao 2021), One Shot Federated Learning (OneshotFL, Guha, Talwalkar, and Smith 2019, in addition to the BCM (Tresp 2000) and the mixture model.

In the case of FedAvg, FedProx, FedBE, and OneshotFL, either SGD with momentum or the Adam optimizer (Kingma and Ba 2015) was used for local optimization (as selected by a grid search). For the rest of the methods, including our own, since they require MCMC sampling, cyclic stochastic gradient Hamiltonian Monte-Carlo (cSGHMC, Zhang et al. 2020) was used. For EP MCMC, the algorithm was computationally intractable for neural network models due to the calculation of the inverse of a covariance matrix over parameters. Thus a diagonal covariance matrix was assumed (which corresponds to the assumption that the local posteriors are approximated by an axis-aligned Gaussian). All methods were run for a single round.

Training Details

For MNIST, Fashion MNIST and EMNIST, the training was run for 25 epochs per client overall (split into 5 rounds for FedAvg and FedPA, while only run in a single round for the

rest). For CIFAR10 and CIFAR100, training was run for 50 epochs per client (split into 10 rounds for multi-round methods). For the regression datasets, a total of 20 epochs are used for all datasets except “air quality”, which used 100 epochs (in both cases, divided into 5 rounds for multi-round methods). The methods involving sampling used a maximum of 6 samples for all experiments.

Hyperparameter Tuning

Hyperparameters were selected based on searching a grid for the best performing settings according to the validation set performance (accuracy for classification, and mean squared error for regression).

The hyperparameters tuned, and their corresponding grids are outlined in Table 3 for both classification and regression. Note that FedPA requires a sampler at each client, and an optimizer at the server. To distinguish where each hyperparameter is used for this algorithm, we therefore label these cases FedPA(C) and FedPA(S) respectively.

In the tables, (D-) β -PredBayes denotes the (distilled) β -predictive Bayes algorithm.

The optimizers used for client training include SGD, SGD with momentum (SGDM), and Adam, while for distillation, we also used Stochastic Weight Averaging (SWA) as a possible optimizer.

The tuned hyperparameter settings for the homogeneous ($h = 0$) classification datasets are in Table 4, while for the heterogeneous setting $h > 0$ they are in Table 5. For the regression datasets, the tuned values are in Table 6.

Note about reading tables: for these tables, if a hyperparameter is repeated more than once, with an algorithm named beside it in brackets, it means the hyperparameter for that algorithm is different. The rest of the algorithms associated with that hyperparameter use the value listed without brackets. For instance, in Table 6, for the “Bike” dataset, the sampler learning rate is listed in the row “Sampler Learning Rate” as $2e-1$, while it is listed separately with the additional specification “FedPA (C)”, as $5e-1$. This means that FedPA uses a sampler learning rate of $5e-1$, while the other sampling algorithms use $2e-1$.

Other hyperparameters not part of the grid search include:

- Batch size: fixed to 100 for all experiments
- Momentum in SGDM: fixed to 0.9 for all experiments
- Model architecture (as described in the main paper)

More algorithm-specific decisions/hyperparameters include:

- FedBE: 10 model samples were drawn from the approximate posterior to use in the ensemble for all experiments (following the experiments in the original paper. This gave a total ensemble size of 16 models = 10 (sampled) + 5 (client models) + 1 (averaged model). By contrast β -predBayes contained an ensemble with 5 models.
- Adaptive FL: The FedYogi server update was used, based on the results from, which suggested that it performed best among their proposed variants. $\beta_1 = 0.9$ and $\beta_2 = 0.99$ were fixed, again, based on the paper.

- One-Shot FL: For the classification case, aggregation is done by averaging the logits of the client models. (This is opposed to averaging and normalizing the probabilities, after the softmax layer).

In tables 3, 4, 5 and 6, LR is used to denote learning rate.

Heterogeneous Classification Dataset Construction

The process for constructing a heterogeneous classification dataset is as follows:

- A parameter $h \in [0, 1]$ is fixed.
- The dataset is sorted by class labels, and split evenly into shards for each client (the “fully heterogeneous shards”).
- A copy of the dataset is made and split such that each shard contains a roughly uniform class distribution for each client (the “homogeneous shards”).
- To form the final shard for a client, a fraction h of each homogeneous shard is replaced with the data from the corresponding fully-heterogeneous shard.

In this way, $h = 0$ corresponds to the homogeneous data setting, while $h = 1.0$ corresponds to a degenerate heterogeneous case (for class distributions in each client).

Regression Models

A note on implementation for regression: the distilled student model is the same architecture as the client models, except in the last layer where it is made to output both a mean and an input-dependent variance. This network is trained to minimize the KL divergence between its output Gaussian distribution, and that of the teacher network.

Additionally, for regression, only the baselines which provide an ensemble are evaluated (since the ensemble is used to predict the variance estimate of the output, $\sigma^2(x^*)$).

Training β

We choose β in β -predBayes to minimize the negative log-likelihood on the server dataset \mathcal{U} . This is done using 10 epochs of gradient descent, using the Adam optimizer with a learning rate of $1e-2$ (for all experiments).

Hyperparameter	Grid Settings		Algorithms Used In
	Classification	Regression	
Optimizer	{SGD, SGDM, Adam}		FedAvg, OneshotFL, FedPA(S), FedProx, AdaptFL, FedBE
Local LR	{1e-1, 1e-2, 1e-3}	{1e-1, 1e-2, 1e-3, 1e-4}	FedAvg, OneshotFL, FedProx, AdaptFL, FedBE
Server LR	{1, 5e-1, 1e-1, 1e-2}		FedPA(S), AdaptFL
Cov. Param (ρ)	{0.4, 0.9, 1.0}		FedPA(C)
Proximal Parameter (λ)	{1, 1e-1, 1e-2, 1e-3}		FedProx
Adaptivity (τ)	{1, 1e-1, 1e-2, 1e-3}		AdaptFL
Sampler LR	{5e-1, 1e-1, 1e-2, 1e-3}	{5e-1, 2e-1, 1e-1, 1e-2, 1e-3}	(D)PB, EPMCMC, FedPA(C)
Maximum Samples	{4, 6, 12}		(D)PB, EPMCMC, FedPA(C)
Temperature	$\{\frac{1}{ \mathcal{D} }\}$	{1, 5e-1, 5e-2, $\frac{1}{ \mathcal{D} }$ }	(D)PB, EPMCMC, FedPA(C)
Sampler Cycles	{5}	{2, 4, 5}	(D)PB, EPMCMC, FedPA(C)
Samples per Cycle	{2}	{1, 2, 3}	(D)PB, EPMCMC, FedPA(C)
Distill Optimizer	{SGDM, Adam, SWA}		D-PB, OneshotFL, FedBE
Distill LR	{1e-2, 5e-3, 1e-4}	{1e-2, 5e-3, 1e-3, 1e-4}	D-PB, OneshotFL, FedBE
Distill Epochs	{100, 50, 20}		D-PB, OneshotFL, FedBE

Table 3: The hyperparameters tuned, their possible values in the grid search, and the algorithms each hyperparameter applies to.

Hyperparameter	Tuned Value				
	MNIST	Fashion MNIST	EMNIST	CIFAR10	CIFAR100
Optimizer	SGDM				
Optimizer (FedProx)	SGDM	Adam	Adam	SGDM	SGDM
Optimizer (AdaptFL)	SGDM	SGDM	SGDM	SGD	SGD
Local LR	1e-1	1e-1	1e-1	1e-2	1e-2
Local LR (FedProx)	1e-1	1e-3	1e-3	1e-2	1e-2
Local LR (AdaptFL)	1e-2	1e-2	1e-2	1e-1	1e-1
Server LR	1	5e-1	1e-1	5e-1	5e-1
Server LR (AdaptFL)	1e-1				
Cov. Param (ρ)	0.4				
Proximal Param (λ)	1e-2	1e-3	1e-3	1e-3	1e-3
Adaptivity (τ)	1e-2				
Sampler LR	5e-1	1e-1	1e-1	1e-1	1e-1
Sampler LR (FedPA(C), EP MCMC)	1e-1				
Maximum Samples	6				
Temperature	$\frac{1}{ D }$				
Sampler Cycles	5				
Samples per cycle	2				
Distill Optimizer	Adam				
Distill LR	1e-4				
Distill Epochs	100				

Table 4: The tuned values of hyperparameters for the classification datasets in the homogeneous case $h = 0$.

Hyperparameter	Tuned Value				
	MNIST	Fashion MNIST	EMNIST	CIFAR10	CIFAR100
Optimizer	SGDM				
Optimizer (FedProx)	Adam	Adam	Adam	SGDM	SGDM
Optimizer (AdaptFL)	SGD				
Local LR	1e-2	1e-2	1e-3	1e-3	1e-3
Local LR (FedProx)	1e-3	1e-3	1e-3	1e-2	1e-2
Local LR (AdaptFL)	1e-1				
Server LR	1	5e-1	1e-1	5e-1	5e-1
Server LR (AdaptFL)	1e-1				
Cov. Param (ρ)	0.4				
Proximal Param (λ)	1e-2	1e-2	1e-2	1e-3	1e-3
Adaptivity (τ)	1e-2				
Sampler LR	1e-1				
Maximum Samples	6				
Temperature	$\frac{1}{ D }$				
Sampler Cycles	5				
Samples per cycle	2				
Distill Optimizer	Adam				
Distill LR	1e-4				
Distill Epochs	100				

Table 5: The tuned values of hyperparameters for the classification datasets, in the heterogeneous case $h > 0$.

Hyperparameter	Tuned Value				
	Air Quality	Bike	Wine Quality	Real Estate	Forest Fire
Optimizer	Adam				
Optimizer (FedPA(S))	SGDM	Adam	SGDM	SGDM	SGDM
Optimizer (FedProx)	SGDM	Adam	Adam	Adam	SGDM
Optimizer (AdaptFL)	SGDM	SGDM	SGD	SGD	SGD
Local LR	1e-2	1e-2	1e-3	1e-2	1e-4
Local LR (FedProx)	1e-2	1e-2	1e-2	1e-1	1e-1
Local LR (AdaptFL)	1e-2	1e-1	1e-1	1e-1	1e-2
Server LR	1e-1	1e-2	5e-1	1	1
Server LR (AdaptFL)	1e-1	1e-1	1	1e-1	1
Cov. Param (ρ)	0.4	1.0	0.9	0.4	0.4
Proximal Param (λ)	1e-2	1e-3	1e-3	1e-2	1e-1
Adaptivity (τ)	1e-2	1e-1	1	1e-2	1e-2
Sampler LR	1e-1	2e-1	2e-1	2e-1	1e-2
Maximum Samples	6	4	4	6	4
Temperature	1	$\frac{1}{ \mathcal{D} }$	5e-2	5e-1	5e-1
Sampler Cycles	5	5	4	5	2
Samples per cycle	1	2	2	2	2
Distill Optimizer	Adam				
Distill LR	1e-3	1e-3	5e-3	5e-3	5e-3
Distill Epochs	100				

Table 6: The tuned values of hyperparameters for the regression datasets.