

Federated Geometric Monte Carlo Clustering to Counter Non-IID Datasets

Federico Lucchetti*, Maria Fernandes[†], Lydia Y. Chen[‡], Jérémie Decouchant[‡], Marcus Völp*

*SnT - University of Luxembourg, [†]WHG - University of Oxford, [‡]Delft University of Technology,
federico.lucchetti@uni.lu, maria.fernandes@well.ox.ac.uk, j.decouchant@tudelft.nl, y.chen-10@tudelft.nl, marcus.voelp@uni.lu

Abstract—Federated learning allows clients to collaboratively train models on datasets that are acquired in different locations and that cannot be exchanged because of their size or regulations. Such collected data is increasingly non-independent and non-identically distributed (non-IID), negatively affecting training accuracy. Previous works tried to mitigate the effects of non-IID datasets on training accuracy, focusing mainly on non-IID labels, however practical datasets often also contain non-IID features. To address both non-IID labels and features, we propose FedGMCC¹, a novel framework where a central server aggregates client models that it can cluster together. FedGMCC clustering relies on a Monte Carlo procedure that samples the output space of client models, infers their position in the weight space on a loss manifold and computes their geometric connection via an affine curve parametrization. FedGMCC aggregates connected models along their path connectivity to produce a richer global model, incorporating knowledge of all connected client models. FedGMCC outperforms FedAvg and FedProx in terms of convergence rates on the EMNIST62 and a genomic sequence classification datasets (by up to +63%). FedGMCC yields an improved accuracy (+4%) on the genomic dataset with respect to CFL, in high non-IID feature space settings and label incongruency.

I. INTRODUCTION

Federated learning (FL) frameworks [1, 2] are commonly used when regulations (such as the GDPR²) or mere data volume prevent data exchanges. Datasets are distributed across the members (clients) of the federation and the global machine learning optimization problem can approximately be split up into smaller sub-problems that are distributed across independently acting clients. Client solutions are aggregated either by a central server or in a distributed manner. However, this approximation only holds in an idealised settings where client datasets are sampled from the same distribution and data samples form independent events. Deviating from independent and identically distributed datasets (IID) poses a challenge to most FL approaches and results in global models that fail to accurately represent the entire aggregated dataset [3–6]. Sources of non-IID-ness can be found in both the label space and the feature space. In the former, for a given feature, different labels might be attributed due to regional differences (e.g., different sentiments in fashion). This is also termed concept shift and leads ultimately to clients with incongruent labelling participating to the training of a

common FL model. In addition, classes may be imbalanced due to one class of labels being over-represented with respect to others. This paper focuses on the latter, i.e., *feature space IID violations*. They appear for example when a particular feature is over-represented in a member’s dataset compared to other datasets. A biobank might sample genes that bias towards a particular local population [7] with the consequence of possibly overlooking regional differences, such as the well known preponderance of a specific gene mutation coding for sickle cell disease in the Sub-Saharan Africa population [8]. For a given label, features might also vary across datasets, such as different handwriting styles for the same alphanumeric character (e.g., 7 with and without bar). Non-IID-ness may lead to feature and label skew [3] and accuracy degradation [4, 6, 9] by leaving the choice between starting from different initialization weights, causing models to diverge, or accepting diversity suppression when starting from the same weights [5].

Besides statistical heterogeneity, FL algorithms need to cope with system heterogeneity, leading to asynchronous model updates or incomplete local training. The first proposed FL attempt, FedAvg [1], aggregates all client models by averaging their weights, leading to the above inaccuracies in the presence of non-IID-ness. Variations [6, 10] of FedAvg have been proposed, with satisfactory results when the non-IID-ness is less pronounced and exclusively in the label space (see Sec. II). Federated Clustering [11–14] tackles this problem by assigning clients to separate clusters and hence limiting the deleterious transfers of negative knowledge between member models that have been trained on distinct datasets. These approaches result in the construction of trained models with reasonable accuracy despite IID violations, however, at the cost of significant communication and/or computational overheads.

In this paper, focusing on non-IID-ness in the feature space, we propose a novel Monte Carlo Clustering approach, called FedGMCC, to counter non-IID-ness without compromising final model accuracy and while maintaining low communication costs and data privacy. We show that in order to train a global FL model in a non-IID setting, the objective function to be minimized has to be conceptually split up into two components: the first encompasses the IID component for which the usual FedAvg solution holds; the second captures the non-IID contribution, which we show can be solved by introducing *interaction* between client models. We derive this

¹Code and genome dataset available at: <https://figshare.com/s/dc2f4280ce012e12f414>

²<https://gdpr.eu/>

interaction by leveraging an observation about the geometry of training loss manifolds [15–18], namely that seemingly different, stationary solutions to the training loss minimization problem, obtained at the individual member sites, are often connected via simple parametric curves where the training loss is approximately flat. The existence of these curves reveals the pair-wise interaction between models needed to solve the non-IID problem, which leads us to establish a criterion for clustering seemingly different solutions and suppressing negative knowledge transfer between non-connected ones. More importantly, we show that drawing from [19], averaging models along the curve parametrization, can produce global models with enhanced accuracy and generalization. As a summary, we make the following contributions.

- 1) We separate the FL objective in an IID and a non-IID components and, based on the geometry of curved training loss manifolds, put forward an Ansatz solution that simultaneously minimizes both components.
- 2) We demonstrate how to construct this solution using a Monte Carlo sampling of the received client model output spaces, and present novel model weight clustering and aggregation rules.
- 3) We evaluate FedGMCC using the EMNIST62 dataset and a genomic dataset with different non-IID-ness. FedGMCC always outperforms FedAvg [1] and FedProx [6]. It outperforms CFL [13] in case of high non-IID-ness.

II. RELATED WORK

Among the myriad of published FL algorithms, most of them rely on clients to upload model weight updates onto a central server that proceeds to aggregate them and to redistribute the result back to clients. The main difference often relies in the way the aggregation rule is applied. FedAvg [1] was first to allow clients to train a global model locally, instead of transferring data. FedAvg returns to a central server only for aggregating all local models by averaging their weights. FedProx [6] inhibits local updates, by adding a proximal term to the loss function, which restrains divergence between local and global model weights. FedBn [10] achieves accurate results on non-IID datasets by batch normalizing the local neural networks’ input layer before aggregation. FedFV [20] detects inconsistent gradient updates and corrects them before the aggregation step. Our approach (FedGMCC) aggregates models by averaging weights along their geometric connection.

Federated clustering groups certain client models into separate clusters to prevent negative knowledge transfer such as the iterative federated clustering algorithm (IFCA) [11] which solves an incongruency that occurs when non-IID datasets exhibit concept shifts by training multiple models on local data and returning the one with the lowest aggregation loss. However, these works primarily focus on label-space non-IID-ness, leading to inaccurate results in the presence of feature-space non-IID-ness. Moreover, training multiple models induce high communication and computation costs.

Merging local models with approximately the same weights (according to cosine similarity, L1 or L2 norm) reduces the

number of models in the ensemble [12, 13]. For example, clustered FL (CFL) [13] merges models if their weights (or gradients) compare well enough based on a specified metric. Again CFL focuses on label non-IID-ness. In contrast, FedGMCC explicitly considers feature-space non-IID-ness. Among the aforementioned FL approaches, CFL is the most reminiscent to our approach, with two differences. First in CFL, the central server applies clustering after multiple FedAvg iterations whereas FedGMCC clusters after each client-server communicating round. Second, CFL’s clustering rule relies on measuring whether the gradients of individual client-weight models are coherent (parallel) via the cosine similarity measure. Our approach verifies that a parallel transport of one client model weights to the other is feasible. Hence, FedGMCC extends CFL’s gradient coherency measure to include intermediary models along the affine connection.

III. PROBLEM DESCRIPTION

We consider K clients that aim to locally minimize a local objective $\mathcal{L}_k = \mathcal{L}(\mathbf{X}_k, \mathbf{w})$. Clients send their respective solutions (client model weights) $\mathbf{w}_k^* = \arg \min_{\mathbf{w}} \mathcal{L}_k$ to a central server to approximate with aggregate \mathbf{w}_t^f a global solution for the loss minimization problem $\mathbf{w}^* = \arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{x}, \mathbf{w})$, which maps the C classes of the compact input space \mathcal{X} into the label space $\mathcal{Y} = [\mathbf{C}]$, where $[\mathbf{C}] = 1, \dots, C$. We consider the neural network as a function $f(\mathbf{w}) : \mathcal{X} \mapsto \mathcal{S}$, parametrized by weights \mathbf{w} , that maps \mathbf{X} to the probability simplex $\mathcal{S} = \{\mathbf{z} | \sum_{i=1}^C z_i = 1, z_i \geq 0, \forall i \in [\mathbf{C}]\}$. We write f_i for the probability of the i -th class and define the population loss $\mathcal{L}(\mathbf{X}_k, \mathbf{w}_k)$ at federation member k as the cross-entropy loss $\mathcal{L}(\mathbf{X}_k, \mathbf{w}) = \frac{1}{N} \sum_{x_i \in \mathbf{X}} \sum_{j=1}^C p(y=j) \mathbb{E}_{\mathbf{x} | y=i} [\log f_j(x_i, \mathbf{w}_k)]$.

In an IID setting, the optimization problem minimizing \mathcal{L}_{IID} can be expressed as K sub-problems that train local models via stochastic gradient descent (SGD), which are then communicated to a server for aggregation after each communication round (indexed by t):

$$\mathcal{L}_{IID} = \sum_{k=1}^K \frac{n_k}{n} \mathcal{L}(\mathbf{X}_k, \mathbf{w}) = \mathbb{E}_{\mathbf{X}_k} (\mathcal{L}(\mathbf{X}_k, \mathbf{w})) \quad (1)$$

FedAvg [1] makes use of the weighted averaging aggregation function $\mathbf{w}_t^f = \sum_k \frac{n_k}{n} \mathbf{w}_{t,k}^*$, where $\frac{n_k}{n}$ is the fraction of the k -th member’s local dataset with respect to the overall dataset size. Unfortunately, in a non-IID setting, this aggregation is not a valid approximation.

We follow [5] in quantifying non-IID-ness with the help of the Earth Mover Distance (EMD) [21] and write $d(\mathcal{D}_1, \dots, \mathcal{D}_K)$ for the EMD of K distributions $\mathcal{D}_1, \dots, \mathcal{D}_K$ obtained by averaging the set of pairwise EMDs (cf. Apx. A).

Weighted averaging fails as an aggregation function under non-IID datasets [5] because: (1) if weights in client models are all initialized identically, the label-space non-IID-ness between client dataset distributions will be the dominant factor and weights will diverge; and (2) if weights are initialized differently, client models risk converging to different solutions

³. The steepness of local loss function pockets in relation to the local flatness of the loss manifold exacerbates the divergence after aggregation and hence leads to sub-optimal global model performance [22], which we address with Geometric Monte Carlo Clustering.

IV. GEOMETRIC MONTE CARLO CLUSTERING

In a non-IID setting, individual client models are set to converge during training towards distinct local minima, i.e., different coordinates in weight space on the training loss manifold (see Fig. 1). In particular, when the non-IID-ness lies exclusively in the feature space between two client datasets $\{\mathbf{X}_1, \mathbf{y}_1\}$ and $\{\mathbf{X}_2, \mathbf{y}_2\}$, we hypothesize the existence of a continuous transformation \mathcal{T} (e.g. pixel rotation, color inversion, shape distortion, etc.) that maps, on average, one subset of features $\mathbf{X}'_1 \subseteq \mathbf{X}_1$ detained by one client to a subset of features $\mathbf{X}'_2 \subseteq \mathbf{X}_2$ of a different client $\mathcal{T}: \mathbf{X}_1 \rightarrow \mathbf{X}_2$. As a consequence, we suppose that the model weights of one client are trained to encode a set of transformed features with respect to another client. Hence we hypothesize the existence of a continuous transformation Γ that maps one subset of weights to another one $\Gamma: \mathbf{w}'_1 \subseteq \mathbf{w}_1 \rightarrow \mathbf{w}'_2 \subseteq \mathbf{w}_2$ and that can alternatively be modeled by a continuous affine connection via a curve parametrization $\gamma_\theta(u)$. This curve connects \mathbf{w}_1 and \mathbf{w}_2 (with $\gamma_\theta(0)=\mathbf{w}_1$ and $\gamma_\theta(1)=\mathbf{w}_2$) ideally on a loss surface where the loss value does not vary. That is, two neural nets parametrized by \mathbf{w}_1 and \mathbf{w}_2 agree on the output space given the same input. Along this curve of invariant loss reside a family of intermediary model weights that can be aggregated to produce a richer global model, incorporating the knowledge of both client models. In contrast, under label-space IID violations, particularly concept shifts, we expect no continuous curve to be found (e.g. between \mathbf{w}_3 and \mathbf{w}_4 on Fig. 1) since they disagree on the classification of a same input feature. In this case, a well constructed clustering rule should separate both models to avoid negative knowledge transfer. This curve finding is at the heart of our novel Federated Geometric Monte Carlo Clustering algorithm (FedGMCC).

A. Prerequisites

Key to FedGMCC is the treatment of non-IID datasets as perturbations of probability distributions \mathcal{D}_k from the ideal IID baseline \mathcal{D} , where $\mathcal{D}_k \mapsto \mathcal{D} + \delta_k$ and $d(\mathcal{D}_k, \mathcal{D}) \neq 0$. Datasets that originally had overlapping feature and label representations have, after perturbation, a component that contributes to their non-IID-ness. Let us denote as $\{\mathbf{X}_k^U, \mathbf{y}_k^U\}$ the labeled dataset drawn from the distribution \mathcal{D}_k^U such that $\forall k \neq l, d(\mathcal{D}_k^U, \mathcal{D}_l^U) = 0$ and as $\{\mathbf{X}_k^\delta, \mathbf{y}_k^\delta\}$ the datasets drawn from \mathcal{D}_k^δ such that $\forall k \neq l, d(\mathcal{D}_k^\delta, \mathcal{D}_l^\delta) \neq 0$, where $\{\mathbf{X}_k^U, \mathbf{y}_k^U\} \cup \{\mathbf{X}_k^\delta, \mathbf{y}_k^\delta\} = \{\mathbf{X}'_k, \mathbf{y}'_k\}$. The training objective is:

$$\begin{aligned} & \sum_k^K \mathcal{L}(\tilde{\mathbf{X}}_k, \mathbf{w}) \\ & \approx \sum_k^K \mathcal{L}(\mathbf{X}_k^U, \mathbf{w}) + \sum_{k < l}^K \mathcal{L}(\mathbf{X}_k^\delta \cup \mathbf{X}_l^\delta, \gamma_{k,l}) \quad (2) \\ & = \mathcal{L}_{IID} + \mathcal{L}_{INT} \end{aligned}$$

³Without a formal proof, a third contribution term to weight divergence could be added in terms of feature-space non-iid

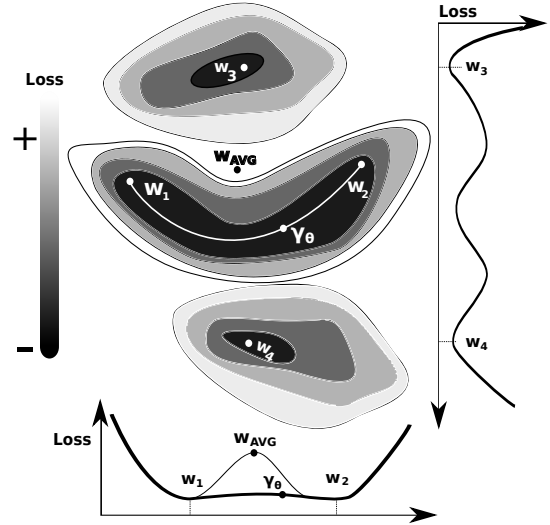


Fig. 1: Schematization of the Geometric Monte Carlo Clustering applied on 4 trained client models where 1 and 2 have a feature space based IID, 3 and 4 have a label space concept shift based non-IID. Curved line connects \mathbf{w}_1 to \mathbf{w}_2 where the loss remains minimally low. No curve is found between \mathbf{w}_3 to \mathbf{w}_4 where the loss crossed the ϵ budget. Note that because of the curved loss surface, the standard FedAvg aggregation $\mathbf{w}_{AVG} = 0.5 \cdot (\mathbf{w}_1 + \mathbf{w}_2)$ would lead to a suboptimal model (high loss).

\mathcal{L}_{IID} is the usual IID loss function (see Eq. 1). The *pairwise interaction loss* \mathcal{L}_{INT} captures the training loss on a pair of non-IID distributions. We use the solutions \mathbf{w}_k to the pure IID objective to approximate the solution to $\min_{\mathbf{w}} \mathcal{L}_{INT}$.

B. Curve Finding

We intuit the solution to the minimization of Eq. 2 to be of the form $\mathbf{w} = (1 - u)\mathbf{w}^f + u\theta_{k,l}$ where $\theta_{k,l}$ introduces a pairwise *interaction* between two client model weights and balances its contribution to the total solution with a coupling constant $u \in [0, 1]$. As we have already hinted in the previous section on the possibility of connecting distinct model weights via a transformation, we propose to model this interaction as an affine connection between weights \mathbf{w}_k and \mathbf{w}_l by a smooth curve $\gamma_\theta(u): [0, 1] \mapsto \mathcal{M}$ on the loss function manifold \mathcal{M} of dimension $|net|$ equal to the number of parameters of the neural net and parametrized by $\theta \in \mathbb{R}^{|net|}$. We need to find the transport parameter θ that leaves the gradient of the second term in Eq. 2 parallel for every $u \in [0, 1]$. This amounts to solving the geodesic equation $\partial_u \nabla_\theta \mathcal{L}(\mathbf{X}_k^\delta \cup \mathbf{X}_l^\delta, \gamma_\theta(u)) = 0$ with boundary conditions $\gamma_\theta(0) = \mathbf{w}_k$ and $\gamma_\theta(1) = \mathbf{w}_l$. Ideally, the central server, which detains the individual client model weights, should execute the curve finding procedure after every local update. However, since the central server does not have access to the datasets in order to explore the loss manifold, we sample a surrogate version of the latter by generating a Monte Carlo type input dataset drawn from a uniform distribution $\mathbf{X}_{MC} \sim U(0, 1)$ to reconstruct the loss from the label-space output $f(\mathbf{X}_{MC}, \mathbf{w}_k)$, using mean-squared error (MSE) as loss function, since we merely compare raw outputs. The curve is

then derived by perturbing the parameter θ in the direction where the loss $\mathcal{L}_{MSE}(\mathbf{w}_l, \gamma_\theta(u))$ does not vary. We rely on a polygonal chain as Ansatz, as it has been proven to lead to optimal curve finding results [15].

$$\gamma_\theta(u) = \begin{cases} 2(u\theta + (0.5 - u)\mathbf{w}_k) & u \in [0, 0.5[\\ 2((u - 0.5)\mathbf{w}_l) + (1 - u)\theta & u \in [0.5, 1] \end{cases} \quad (3)$$

Following [15], we simplify the loss to be minimized to the expectation of $\mathcal{L}_{MSE}(\mathbf{w}_l, \gamma_\theta(u))$ with respect to uniform distribution on the curve on $u \in [0, 1]$, making the loss computationally tractable.

$$\nabla_\theta \int_0^1 \mathcal{L}_{MSE}(\mathbf{w}_l, \gamma_\theta(u)) du = \frac{\nabla_\theta \mathbb{E}_{u \sim U(0,1)} \mathcal{L}_{MSE}(\mathbf{w}_l, \gamma_\theta(u))}{\nabla_\theta \mathbb{E}_{u \sim U(0,1)} \mathcal{L}_{MSE}(\mathbf{w}_l, \gamma_\theta(u))} \quad (4)$$

C. Model Weights Clustering and Aggregation

Garipov et. al [15] showed that two models initialized differently and trained on the same dataset can be connected via a simple curve. Two models trained independently on distributed datasets with only partially overlapping features are also expected to converge to different stationary solutions. Nevertheless, we expect these seemingly different solutions to be interlinked via an affine connection in weight space i.e a smooth curve can be constructed that along a approximately flat loss manifold. Hence, seemingly different solution form a family of equivalent solutions to the loss minimization problem which leads us to put forward the following clustering rule.

Proposition 1.1 $\mathbf{X}_{MC} \sim U(0, 1)$, two models with weights \mathbf{w}_k and \mathbf{w}_l belong to the same cluster \mathcal{S} if there exists a $\theta \in \mathbb{R}^{|\text{net}|}$ that parametrizes the curve $\gamma_\theta : u \in [0, 1] \mapsto \mathbb{R}^{|\text{net}|}$ and $\epsilon \geq 0$ such that

$$\partial_u \mathbb{E}_{x \in \mathbf{X}_{MC}} \|f(x, \mathbf{w}_l) - f(x, \gamma_\theta(u))\|^2 < \epsilon \quad (5)$$

Because $\gamma_\theta(u)$ generates a family of intermediary model weights, all approximate solutions to the optimization problem (see Eq. 2), their sum must also be a solution. [19] showed that averaging intermediary model weights on the constructed curve $\gamma_\theta(t)$ along t can lead to an aggregated model with improved generalization.

Proposition 1.2 Model weights \mathbf{w}_j belonging to the same cluster \mathcal{S}_j are aggregated with

$$\mathbf{w}_j = \sum_{k,l \in \mathcal{S}_j} \mathbb{E}_{u \in [0,1]} \gamma_{\theta_{k,l}}(u) \quad (6)$$

D. Algorithm

FedGMCC is described in Alg. 1 and schematized in Fig. 1. The server samples the output space of each received client model by generating a Monte Carlo input dataset and tests via Prop. 1.1 whether two models can be grouped together. This procedure can be seen as a disjoint-set query in which model weights that could not be connected to any other model weights are kept as singleton sets (e.g., $\{\mathbf{w}_1, \mathbf{w}_2\}$, $\{\mathbf{w}_3\}$, $\{\mathbf{w}_4\}$ on Fig. 1). With M disjoint sets, the weights that belong to

Algorithm 1 Geometric Monte Carlo Clustering

```

1: Parameters:  $\epsilon > 0$ ;  $n$ : Monte Carlo sample size;  $K$ : number of client
   models;  $M$ : number of clusters;  $\eta$ : learning rate
2: Initialization: Random input dataset  $\mathbf{X} \leftarrow \text{Uniform}(0, 1, n)$ 
3: Server: ModelClustering( $\mathbf{w}_1, \dots, \mathbf{w}_K$ ):
4:   clusters  $\mathcal{S} \leftarrow \{\}$ 
5:   for  $k \in 1, \dots, K, l \in k + 1, \dots, K$  do
6:     if ( $\{\mathbf{w}_k, \mathbf{w}_l\} \in \mathcal{S}_{j \leq k}$ ) continue
7:     else
8:        $\theta_{k,l} \leftarrow \text{Update}(\mathbf{w}_k, \mathbf{w}_l)$ 
9:       if ( $\theta_{k,l}$ )  $\mathcal{S}_j \leftarrow \mathcal{S}_j \cup \{\mathbf{w}_k, \theta_{k,l}\}$ 
10:       $\mathbf{w}_1, \dots, \mathbf{w}_M \leftarrow \text{Aggregation}(\mathcal{S})$ 
11:      Distribute  $\mathbf{w}_1, \dots, \mathbf{w}_M$ 
12:    Update( $\mathbf{w}_1, \mathbf{w}_2$ ):
13:       $\theta \leftarrow \mathbf{w}_1 + \mathbf{w}_2$ 
14:      for  $u \in U(0, 1), x \in \mathbf{X}_{MC}$  do
15:         $\mathcal{L}(x, \mathbf{w}_2, \gamma_\theta(u)) \leftarrow \|f(x, \mathbf{w}_2) - f(x, \gamma_\theta(u))\|^2$ 
16:         $\theta_{u+1} \leftarrow \theta_u - \eta \nabla_\theta \mathcal{L}(x, \mathbf{w}_2, \gamma_\theta(u))$ 
17:        if ( $\max_u \partial_u \mathbb{E}_{x \in \mathbf{X}_{MC}} \mathcal{L} \leq \epsilon$ ) return  $\theta$ 
18:      Aggregation( $\mathcal{S}_{1, \dots, M}$ ):
19:        for  $\mathcal{S}_j \in \mathcal{S}$  do  $\{\mathbf{w}_j \leftarrow \sum_{k,l \in \mathcal{S}_j} \mathbb{E}_{u \in U(0,1)} \gamma_{\theta_{k,l}}(u)\}$ 
20:        return  $\mathbf{w}_1, \dots, \mathbf{w}_M$ 
21:      Distribute( $\mathbf{w}_1, \dots, \mathbf{w}_M$ ):
22:        for  $\mathcal{S}_j \in \mathcal{S}, k \in 1, \dots, K$  do
23:          if ( $k \in \mathcal{S}_j$ ) send  $\mathbf{w}_j$  to client  $k$ 

```

a same cluster \mathcal{S}_j are aggregated by averaging all pairwise interactions leading to M clustered global solutions to the non-IID problem (see Proposition 1.2). \mathbf{w}_j is then sent to the clients that contributed to it.

V. EVALUATION

To evaluate our approach, we show that real-life data is in fact non-IID in the feature space and compare the accuracy of our approach — Federated Geometric Monte-Carlo Clustering (FedGMCC) — against state-of-the art federated learning approaches (FedAVG [1], FedProx [6], CFL [13] and standard SGD). We leverage the well known image dataset EMNIST62 [23] and two sequential genomic datasets (SENSG-A and SENSG-R) we created (see Appx. A). EMNIST62 contains 814255 handwritten characters, labeled as 62 unbalanced classes in 28x28 pixel format. The genomic dataset contains reads (i.e., sequences of the bases A, T, C, G) of size 150 and every position is labelled according to whether or not it is part of a genomic variation. Such labelling is important, for example, to filter out and protect private information in the genomic information processing pipeline [24]. We introduce a concept shift in SENSG-R by randomly flipping the label of the most recurrent genomic features. EMNIST62 and SENSG-A were distributed among K clients following the procedure detailed in Appx. A. This partitioning introduces non-IID-ness in the feature space but also in the label space because certain classes might be over-represented in some client datasets compared to others. The genomic dataset SENSG-R is naturally partitioned using the populations of individuals (Asian, European, African). In addition, we partition both sets artificially into 10 client datasets to consider a wider range of non-IID-ness. Models were trained on an AMD Ryzen7 3700x system with 8 3.6 GHz cores, a NVIDIA Geforce RTX 3090 GPU with 10496 CUDA cores and 24 GB of GDDR6X

memory. We used the binary-cross entropy loss function for the training of all classifiers and Tensorflow 2.6 [25].

A. Network Architecture and Baselines

We construct two neural network architectures for the benchmarks: for EMNIST62, we use 2 stacked CNN layers, activated by ReLu, and followed by a fully connected neural network. The genomic datasets are classified with a network comprised of 2 stacked bidirectional LSTMs that feed into a densely connected neural network. We compare our approach FedGMCC against three federated learning methods: FedAvg [1], FedProx [6] and CFL [13]. We assume a federation of $K = 10$ members for all experiments, each solving its local optimization problem using SGD. FedAvg performs weighted average aggregation in the central server after each local training round, which is performed on the member datasets. FedProx [6] corrects FedAvg’s loss function by the proximal term $\mu \|w - w_t\|^2$, where μ is a positive constant, w the weights of the most recent global model and w_t the weights of the local model at training step t . We implement CFL by following the procedure laid out in [13]. With FedAvg the central server receives, aggregates and distributes client updates until the average of received client gradients decreases below $\epsilon_1 = 0.2$. FedGMCC applies the clustering, aggregation and distribution scheme laid out in Alg. 1. In the first round, the value of ϵ was set to the median value of the $\frac{1}{2}K(K-1)$ training losses associated with the curve finding procedure. This value is subject to a 5 percentile increase if the number of clusters was higher than 1 and a 5 percentile decreases otherwise. The learning rate η in the curve finding procedure was set to 0.1. In addition, we train two central models ($cSDG_0$ and $cSDG_1$) on the combined genomic dataset using standard SGD as baselines. Training of two models is necessary due to the concept shift we introduce in SENS-G-A. As usual, we separate the whole dataset into disjoint training (80%) and validation sets (20%). Hyperparameters have been fine-tuned for every training algorithm to give best possible results in terms of convergence rate. The mini-batch size was set to 64, and learning rates were set to 0.001. Local epoch numbers were respectively set to 10 and 5 for the EMNIST62 and SENS-G-A datasets.

B. Results

Curve Fitting between Local Models: Fig. 2 shows the average training loss for our curve-fitting approach for different parameters and setups. The top graph compares a naive linear curve $\gamma(t) = u w_1 + (1-u) w_2$ with a polygonal one-bend chain curve (see 3). Using the latter we were able to find a region of low loss and hence prove the existence of multiple simple connections between pairs of client models, i.e., pairs of client models belonging to the same cluster. Note that by setting $u = 0.5$, the linear parametrization reduces to the classical FedAvg aggregated model $w_1 + w_2$ situated in a region of higher loss with respect to the chain curve loss $\gamma_\theta(u = 0.5)$. When clients train local models with varying local epoch numbers (middle left) or EMD values

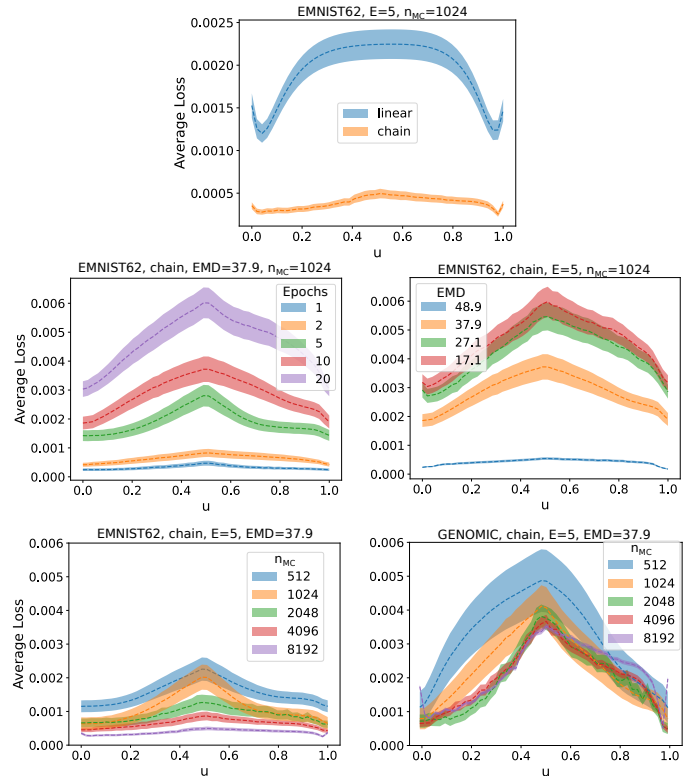


Fig. 2: The 10 client averaged train loss as along the curve $\gamma_\theta(u)$ connecting two client models as a function of multiple experimental setting (EMD, local epochs, datasets) and curve parameters (type, Monte Carlo sample size n_{MC}). Dispersion represents the standard deviation scaled down by a factor of 10.

(middle right) the average loss is directly affected. When the number of local epochs is decreased or the data distribution setting tends towards the ideal IID setting, the GMCC losses flatten out and the inter-client variability decreases. This is an important finding that could previously only be obtained when computing loss surface manifold with real datasets. Our approach achieves this on a generated dataset using Monte Carlo sampling. More interestingly, increasing the sample size of the GMCC input dataset makes it easier to find connections between client model weights with low training losses. For EMNIST62, where no concept shift was present (bottom left), finding an optimal curve parametrization is only a matter of increasing computational costs (due to increased GMCC sample sizes). On the other hand, for the genomic datasets, because the concept-shift introduced an incongruency in the output space between two models, no amount of GMCC sample size will be able to connect models hence leading to the splitting of client models into separate clusters. We validate our approach, where the server substitutes the client datasets by a surrogate Monte Carlo type generated dataset.

Loss: Fig. 3 reports the loss of the FL models for varying EMD values and for the EMNIST62 and genomic datasets. Loss is computed on the validation set after every communication round. For CFL and FedGMCC, where the aggregation can lead to multiple models, the average loss is shown. We

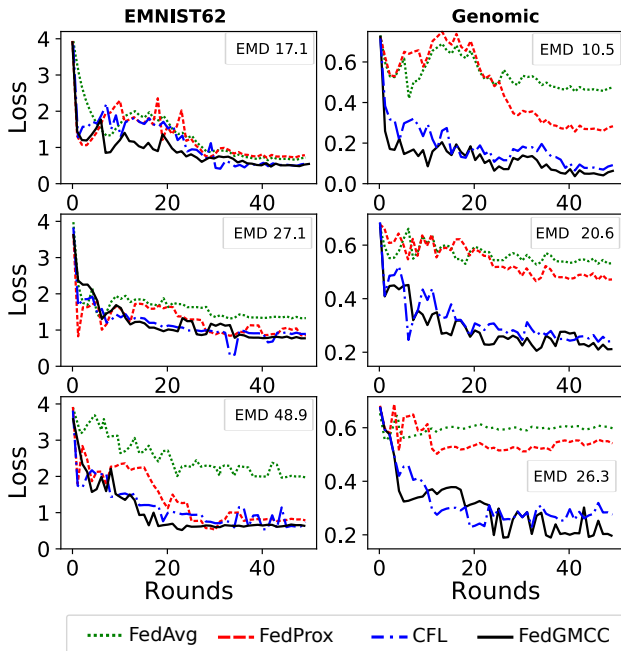


Fig. 3: Validation loss for FedAvg, FedProx, CFL and FedGMCC for three different EMD values.

highlight the increasing difficulty for FL models to converge when the EMD increases. As expected, FedAvg shows the worst results for both datasets. FedProx is able to deal with high EMD values in the EMNIST62 dataset but does not converge to a lower loss value at high EMD compared to the FL clustering algorithms (CFL and FedGMCC). The latter are able to cluster dissimilar models and avoid negative transfer of knowledge between them, which allows them to obtain the two aggregated final models with optimal loss.

Accuracy: For the EMNIST62 dataset the FL algorithms achieved a similar performance with low EMD values than the one obtained in the centralized setting (85%) (see Table I). However, their accuracy degrades quickly when the EMD increases, with values for FedAvg that dropped from 0.81 at EMD 17.1 to 0.43 at EMD 48.9. This drop was less pronounced for FedProx, CFL and FedGMCC. The two latter maintained a 0.79 accuracy at high EMD. The situation was different for the genomic datasets. The centralized SGD models $cSGD_0$ and $cSGD_1$ achieved a 0.88 accuracy. CFL and FedGMCC attained comparable 0.85 accuracy, while FedProx and FedAvg respectively obtained 0.71 and 0.75 at low EMD. As expected, FedAvg’s performance dramatically degraded at high EMD to 0.18. The drop in accuracy was less severe for FedProx and stayed around 0.79. FedGMCC led to two global models, each of which maintained an accuracy of 0.85 at EMD 17.1 and 0.84 at higher EMD values. This high performance was also achieved by CFL but FedGMCC yielded a better accuracy at high EMD. This is not surprising because of CFL’s and FedGMCC’s aggregation rules, which enable them to create multiple personalized models to accommodate a certain degree of IID-ness in particular the incongruency due to the concept

TABLE I: FL algorithms accuracies in different EMD settings. Pair of values for CFL and FedGMCC indicate accuracies associated to the final weight clusters.

EMNIST62				
EMD	FedAvg	FedProx	CFL	FedGMCC
17.1	0.81	0.83	0.83	0.83
27.1	0.76	0.83	0.83	0.83
48.9	0.43	0.71	0.79	0.79
GENOMIC				
EMD	FedAvg	FedProx	CFL	FedGMCC
6.2 (real)	0.88	0.89	0.91	0.93
10.5	0.88	0.91	0.93 0.92	0.93 0.93
20.6	0.56	0.70	0.80 0.79	0.84 0.83
26.3	0.21	0.65	0.80 0.79	0.84 0.83

shift injected in the genomic dataset.

VI. CONCLUSION

We presented FedGMCC, a federated learning framework that consists of novel clustering rules and a new aggregation procedure. Substituting real datasets by a surrogate Monte Carlo dataset, we show how the curve finding procedure can reveal the geometric connection between congruent models serving as a clustering rule of client model weights. Furthermore, FedGMCC aggregation rule averages all the intermediary model weights on the curve parametrization, leading to better generalization and extending the classical FedAvg aggregation rule to weight spaces of curves training loss manifold.

FedGMCC outperformed other FL training algorithms on the EMNIST62 and genomic sequence sensitivity classification tasks where we controlled the non-IID-ness using an artificial partitioning technique and the EMD measure. In high non-IID setting FedGMCC yielded convergence rates respectively 36% and 8% faster than those of FedAvg and FedProx on EMNIST62. FedGMCC outperformed FedAvg, FedProx and CFL on the genomic datasets by 63%, 19% and 4 %.

REFERENCES

- [1] H. Brendan McMahan, Eider Moore, Daniel Ramage, and Blaise Agüera y Arcas. “Federated Learning of Deep Networks using Model Averaging”. In: *CoRR* abs/1602.05629 (2016). arXiv: 1602.05629.
- [2] He Yang. “H-FL: A Hierarchical Communication-Efficient and Privacy-Protected Architecture for Federated Learning”. In: *Proc. of 13th Int. Conf. on Artificial Intelligence, IJCAI, Montreal, Canada, 19-27 August*. 2021, pp. 479–485.
- [3] Peter Kairouz, H Brendan McMahan, Brendan Avent, et al. “Advances and open problems in federated learning”. In: *arXiv:1912.04977* (2019).
- [4] Kevin Hsieh, Amar Phanishayee, Onur Mutlu, and Phillip Gibbons. “The non-iid data quagmire of decentralized machine learning”. In: *Int. Conf. on Machine Learning*. PMLR. 2020, pp. 4387–4398.
- [5] Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. “Federated learning with non-iid data”. In: *arXiv preprint arXiv:1806.00582* (2018).
- [6] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. “Federated Optimization in Heterogeneous Networks”. In: *Proc. of Machine Learning and Systems (MLSys), Austin, TX, USA, March 2-4*. 2020.

- [7] Jeroen GJ van Rooij, Mila Jhamai, Pascal P Arp, et al. “Population-specific genetic variation in large sequencing data sets: why more data is still better”. In: *Eu. J. of Human Genetics* 25.10 (2017), pp. 1173–1175.
- [8] Thomas N Williams. “Sickle cell disease in sub-Saharan Africa”. In: *Hematology/Oncology Clinics* 30.2 (2016), pp. 343–358.
- [9] Qinbin Li, Yiqun Diao, Quan Chen, and Bingsheng He. “Federated learning on non-iid data silos: An experimental study”. In: *arXiv preprint arXiv:2102.02079* (2021).
- [10] Xiaoxiao Li, Meirui Jiang, Xiaofei Zhang, Michael Kamp, and Qi Dou. “Fedbn: Federated learning on non-iid features via local batch normalization”. In: *arXiv preprint arXiv:2102.07623* (2021).
- [11] Avishek Ghosh, Jichan Chung, Dong Yin, and Kannan Ramchandran. “An efficient framework for clustered federated learning”. In: *34th Conf. on Neural Information Processing Systems (NeurIPS)*. Canada, 2020.
- [12] Christopher Briggs, Zhong Fan, and Peter Andras. “Federated learning with hierarchical clustering of local updates to improve training on non-IID data”. In: *IEEE Int. Conf. on Neural Networks (IJCNN 2020)*, pp. 1–9.
- [13] Felix Sattler, Klaus-Robert Müller, and Wojciech Samek. “Clustered federated learning: Model-agnostic distributed multitask optimization under privacy constraints”. In: *IEEE T. on neural net. and learning sys.* (2020).
- [14] Kavya Koppurapu and Eric Lin. “Fedfmc: Sequential efficient federated learning on non-iid data”. In: *arXiv preprint arXiv:2006.10937* (2020).
- [15] Timur Garipov, Pavel Izmailov, Dmitrii Podoprikin, Dmitry Vetrov, and Andrew Gordon Wilson. “Loss surfaces, mode connectivity, and fast ensembling of dnns”. In: *Proc. of the 32nd Int. Conf. on Neural Information Processing Systems*. 2018, pp. 8803–8812.
- [16] C. Daniel Freeman and Joan Bruna. “Topology and Geometry of Half-Rectified Network Optimization”. In: *5th Int. Conf. on Learning Representations, (ICLR), Toulon, France, April 24-26, 2017*.
- [17] Seyed Iman Mirzadeh, Mehrdad Farajtabar, Dilan Gorur, Razvan Pascanu, and Hassan Ghasemzadeh. “Linear mode connectivity in multitask and continual learning”. In: *arXiv preprint arXiv:2010.04495* (2020).
- [18] Anna Choromanska, Mikael Henaff, Michael Mathieu, Gérard Ben Arous, and Yann LeCun. “The loss surfaces of multilayer networks”. In: *Artificial intelligence and statistics*. PMLR. 2015, pp. 192–204.
- [19] Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. “Averaging weights leads to wider optima and better generalization”. In: *Conf. on Uncertainty in Artificial Intelligence, Monterey, CA, USA* (2018).
- [20] Zheng Wang, Xiaoliang Fan, Jianzhong Qi, Chenglu Wen, Cheng Wang, and Rongshan Yu. “Federated Learning with Fair Averaging”. In: *Proc. of 13th Int. Conf. on Artificial Intelligence, IJCAI, Montreal, Canada, 19-27 August, 2021*, pp. 1615–1623.
- [21] Yossi Rubner, Carlo Tomasi, and Leonidas J Guibas. “A metric for distributions with applications to image databases”. In: *6th Int. Conf. on Comp. Vision (IEEE Cat. No. 98CH36271)*. 1998, pp. 59–66.
- [22] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. “On large-batch training for deep learning: Generalization gap and sharp minima”. In: *arXiv preprint arXiv:1609.04836* (2016).
- [23] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. “EMNIST: Extending MNIST to handwritten letters”. In: *Int. Conf. on Neural Networks (IJCNN)* (2017).
- [24] Jérémie Decouchant, Maria Fernandes, Marcus Völp, Francisco M Couto, and Paulo Verissimo. “Accurate filtering of privacy-sensitive information in raw genomic data”. In: *J. of Biomed. Informatics* 82 (2018), pp. 1–12.
- [25] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. “Tensorflow: A system for large-scale machine learning”. In: *12th USENIX Symp. on Operating Systems Design and Implementation (OSDI 16)*. 2016, pp. 265–283.

APPENDIX

Earth Mover Distance (EMD) is defined as a distance metric between two distributions \mathcal{A} and \mathcal{B} . It computes the minimal distance for mapping all clusters of distribution \mathcal{A} (set of suppliers) to any cluster of distribution \mathcal{B} (set of consumers). The distributions are thereby characterized by signatures $\mathcal{A} = \{(p_1, w_{p_1}), \dots, (p_m, w_{p_m})\}$ and likewise $\mathcal{B} = \{(q_1, w_{q_1}), \dots, (q_n, w_{q_n})\}$, where clusters are represented as bin centroids p_i with weight w_{p_i} (and q_j with weight w_{q_j} , respectively). The overall cost for transferring all clusters from \mathcal{A} to \mathcal{B} is:

$$C = \sum_{i \in \mathcal{A}} \sum_{j \in \mathcal{B}} c_{ij} f_{ij} \quad (7)$$

where c_{ij} is the ground distance between the supports p_i and q_j and f_{ij} is the flow between p_i and q_j that needs to be minimized under the constraints:

- 1) $f_{ij} \geq 0$ (unidirectional flow),
- 2) $\sum_{(p_i, w_i) \in \mathcal{A}} f_{ij} \leq w_i$ (limited consumer storage),
- 3) $\sum_{(q_j, w_j) \in \mathcal{B}} f_{ij} \leq w_j$ (limited supply), and
- 4) $\sum_{(p_i, w_i) \in \mathcal{A}} \sum_{(q_j, w_j) \in \mathcal{B}} f_{ij} = \sum_{(p_i, w_i) \in \mathcal{A}} w_i = \sum_{(q_j, w_j) \in \mathcal{B}} w_j$ (transfer all).

With optimal flow $f^* = \arg \min_f (C)$, EMD for a pair of distributions follows as:

$$d(\mathcal{A}, \mathcal{B}) = \frac{\sum_{i \in \mathcal{A}} \sum_{j \in \mathcal{B}} c_{ij} f_{ij}^*}{\sum_{i \in \mathcal{A}} \sum_{j \in \mathcal{B}} f_{ij}^*} \quad (8)$$

We extend $d(\mathcal{A}, \mathcal{B})$ to a set of K distributions $\{\mathcal{D}_1, \dots, \mathcal{D}_K\}$ by averaging the set of pairwise EMDs between the k -th distribution \mathcal{D}_k and the distribution over the whole population \mathcal{D} \mathcal{D}_k , using

$$d(\mathcal{D}_1, \dots, \mathcal{D}_K) = \frac{1}{N} \sum_{k=1}^K n_k d(\mathcal{D}_k, \mathcal{D}) \quad (9)$$

where n_k is the sample size of the dataset generated by distribution \mathcal{D}_k , N is the total population size and $\mathcal{D} = \frac{1}{N} \sum_{k=1}^K n_k \mathcal{D}_k$. In the context where \mathcal{D}_k models a labeled dataset with C classes, the joint probability distribution is factorized in terms of its conditional and marginal probability distribution $\mathcal{D}_k = p_k(\mathbf{X}, \mathbf{y}) = \{p_k(y = i, \mathbf{X}) \cdot p_k(\mathbf{X})\}$ for $i = 1, \dots, C$, i.e, the EMD is computed over all the classes for a given feature \mathbf{X} .

To accelerate EMD computation we reduce datasets to their essentials by training autoencoders using a SGD optimizer and mean square error as reconstruction loss function. Table II

TABLE II: Parameters and network configuration for autoencoders to accelerate EMD computation. Training progressed at a rate of 0.01 for EMNIST62 and with 0.001 for the different genomic datasets. We trained in batches of size 128.

	Encoder Layers	dim(features)	Decoder layers	Learning rate	Reconstruction loss
EMNIST62	4 x CNN	7x7x2	4 x TransCNN	0.01	$\leq 0.0001\%$
SENSG-*	2 x CNN + 2 x LSTM	30x1	2 x LSTM + 2xTransCNN	0.001	$\leq 0.0001\%$

shows the parameters for the autoencoders for the two benchmarks we used to evaluate our approach.

A. Image Datasets

EMNIST62 ⁴ is an image dataset for simulating non-IID image classification [23]. It comprises a set of 814255 hand-written alpha-numeric characters, labeled as 62 unbalanced classes, formatted in 28x28 pixel images.

B. Genomic Datasets

We used two different genomic datasets: SENSG-R and SENSG-S.

SENSG-R is composed of 7 500 000 reads from four randomly selected genomes from each of the three major populations represented in the 1000 Genomes Project (1000GP) ⁵: African, European and Asian(see Table III). With this we intend to resemble the natural representation one would obtain when sampling in regions where these populations are dominant.

SENSG-S is composed of one million reads generated from twenty individual genomes (10^5 reads for each genome) and the randomly selected individuals are the following: HG00096, HG00097, HG00099, HG00100, HG00101, HG00102, HG00103, HG00105, HG00106, HG00107, NA21128, NA21129, NA21130, NA21133, NA21135, NA21137, NA21141, NA21142, NA21143, and NA21144.

TABLE III: Genomes used per population dataset.

ASIA	AFRI	EURO
HG00543	HG02703	HG00315
HG00559	HG02769	HG00327
HG00566	HG02715	HG00334
HG00578	HG02771	HG00339
HG00581	HG02722	HG00341
HG00580	HG02676	HG00346
HG00593	HG02808	HG00353
HG00598	HG02810	HG00358
HG00592	HG02614	HG00360
HG00613	HG02839	HG00365

For generating these datasets, we follow three steps:

Step 1 - Genomic sequence generation: We compiled the chromosome 1 sequence of each individual selected by

combining the human reference genome GRCh37 and the individual’s variants in the 1000 GP. Second, we randomly select positions in the sequence to use as seed for obtaining 150 character sequences (reads), comprised of the nucleotides A, T, G, and C. Reads constitute the first digitized information obtained from next generation sequencing machines, which

⁴https://colab.research.google.com/drive/1r-c6UTkJEQx3Pi-HI9q_MoveIF_0h03M

⁵<https://www.internationalgenome.org/>

are subsequently processed to extract variations (i.e., what distinguishes us one from another and what might carry sensitive information, like, disease prepositions). Therefore, genomic data is a practical example where data is generated in multiple geo-distributed locations, e.g., hospitals in different continents, and it must not be shared among different locations for privacy reasons.

Step 2 – Labeling: We labeled each nucleotide by deeming it as sensitive if a genomic variation reported in the 1000 GP or insensitive otherwise, respectively, 1 or 0. For SENSG-S only, we simulate a concept shift by flipping some labels manually in order to obtain non-IID-ness.

Step 3 – Reads encoding: Next, we used word2vect encoding to convert the genomic sequence in a vector. After, this step the reads are ready to be used for the training. In this step we consider a vocabulary (word size) of 5 letters and we generated all the consecutive 5 nucleotides sequences of each read.

C. Artificial Partitioning

SENSG-A and EMNIST62 are partitioned artificially into K subsets using our iterative clustering and distribution algorithm, which we describe in the following. We first apply k -mean clustering to divide the dataset into K subsets with maximal EMD. Then, computing the center of gravity (COG) of all subsets, we randomly select elements from the farthest subset D_i and assign them to other subsets D_j . This reduces the inter-subset distance $d(D_i, D_j)$ and hence the overall EMD. We retain the current partition if pairwise EMD values ($d(D_i, D_j)$) are normally distributed (according to the Shapiro normality test). We repeat this process until a we found a partition with a total EMD value that is low enough.