



Machine Learning and Optimization Laboratory

Incentivized, Privacy-preserving, and Personalized
Distributed Machine Learning Framework for Medical Data

Semester project — 18/06/2020

Felix Grimberg
Author

Mary-Anne Hartley
Supervisor

Martin Jaggi
Co-supervisor

Sai Praneeth Reddy
Karimireddy
Co-supervisor

Andres Colubri
Co-supervisor

Abstract

The federated learning setting is prone to suffering from non-identically distributed data across participating agents. This gives rise to the task of model personalization, where agents collaborate to train several different machine learning models instead of training only one global model. The aim of model personalization is to minimize the sum of the generalization error incurred from training on small data sets, and the transfer error incurred from applying a globally-trained model to a specific local distribution. In this report, two novel approaches to personalized cross-silo federated learning are introduced and discussed from a theoretical perspective: the adapted Ndoye factor, and the Weight Erosion aggregation scheme. The latter is implemented and compared to two baseline aggregation schemes in two case studies: training a diagnostic model on real-world medical data, and predicting the survival of passengers on the publicly available Titanic data set. The models trained using the Weight Erosion aggregation scheme are compared to those trained using the baseline aggregation schemes, both in terms of their classification accuracy on the local test set and in terms of the learned parameters. We demonstrate that the novel Weight Erosion scheme can outperform both baseline aggregation schemes for some specific tasks.

Contents

1	Aim	2
1.1	Necessary features	2
1.2	Performance features	2
2	Objectives	2
2.1	Landscape analysis	3
2.2	Model personalization and data selection	3
2.3	Application case study	4
3	Landscape analysis	4
3.1	Privacy in ML	4
3.2	Model personalization	4
3.3	Incentives	5
3.4	Reproducible data science	5
4	Model personalization and data selection	6
4.1	Setting and Objective	6
4.2	Summarizing evaluation of methods	7
4.3	The adapted Ndoeye factor	8
4.4	Gradient-based Weight Erosion	9
5	Application case study	12
5.1	Implementation	12
5.2	Case study	13
5.3	Results	13
6	Conclusion	21

1 Aim

This project is part of a larger endeavour to create a platform for medical professionals around the globe, and especially in low-resource settings, to train machine learning (ML) models collaboratively. Main requirements to achieve this goal are identified and separated into necessary features and performance features.

1.1 Necessary features

Data privacy. First and foremost, the proposed platform must safeguard each participant from the leakage and re-identification of individual patients stored within their data set. This implies maintaining data at its original location and limiting transfers, which motivates the use of so-called cross-silo federated learning (FL) techniques [10].

Incentives and intellectual property. Secondly, the platform must present incentives for the collection of high-quality health data and protect the intellectual property of their owner.

Resilience. Moreover, it must be designed to work with incomplete and non-identically distributed data, as well as with arbitrary connection/disconnection patterns from participants.

Personalization. Finally, it is important that the selection of training data, features, and model architecture, is interpretable to the users (i.e., medical professionals) and subjected to their judgement. To this end, users must be given the tools to ensure that the trained ML model is adapted to the particularities of their local population.

1.2 Performance features

Data analysis. For better model personalization, the proposed platform should facilitate the selection of appropriate training data, e.g., by providing and visualizing measures of similarity between data sets.

Feedback. Additionally, it should give its users feedback on the quality of their data collection. For instance, it could point out physiological or epidemiological anomalies (e.g., breathing rates being consistently reported higher than usual). Likewise, it could suggest which additional features should be recorded to efficiently maximize future model performance and thus increase their collaborative scope.

Environmental impact. The framework should also be optimised to quantify, minimise and possibly compensate for the environmental impact associated with training each ML model.

Communication efficiency. The communication effort exerted by each participating device should be kept as low as possible, to facilitate participation in low-infrastructure settings. This has the additional benefit of reducing the environmental impact of each training run.

Robustness. Finally, the design should be robust to a small fraction of non-collaborative participants. Different robustness models exist for various types of non-collaboration, ranging from the inadvertent provision and use of false data to fully adversarial behaviour.

2 Objectives

Within the scope of this project, our main focus will be on developing and evaluating methods for model personalization. Throughout the project, we will be mindful of the other requirements identified in Section 1 and seek methods that are compatible with, or even beneficial for, as many of them as possible. We will strive to achieve these objectives in a fully reproducible manner, while maintaining the confidentiality of the provided data sets.

2.1 Landscape analysis

We will begin by researching existing methods for privacy-preserving distributed machine learning and model personalization. During this landscape analysis, we will also explore available tools for reproducible data science, as well as existing incentive schemes, and investigate how these tools can be leveraged to efficiently realise the other objectives listed herein. Our findings will be documented in Section 3 of this report.

2.2 Model personalization and data selection

Often, each agent in a FL setting collects data on a different population. For instance, a hospital in Freetown and an Ebola treatment center in rural Sierra Leone would see fairly different patients in general, which would again differ from the patients admitted in Ebola treatment facilities in the Democratic Republic of the Congo. While it is also highly useful to make global inference across populations, clinicians in the Freetown hospital may be more interested in being able to accurately predict outcomes for new patients admitted *to their hospital*, than in predicting outcomes for new patients admitted *to any facility in general*. They would still need to use each other’s data collaboratively, simply because the amount of data collected by each user is limited.

We will propose a variety of novel ways to help users in optimising the selected model for their local population. Some of the proposed methods will involve ranking the data sets of other users, exploiting various notions of similarity to the local data set. Novel ideas that we will pursue are:

- **Ndoye:** Computing a (potentially adapted) Ndoye factor for each available user before training [16] (cf. Section 4.3).
- **Influence:** Using a subset of the local data set as a test set, and selecting or weighting other users according to how much they contribute to reducing the test loss. As an additional advantage, this would guarantee that the evaluation metrics of the resulting model measure its ability to make predictions on the local population (as opposed to the global population or some aggregated subset thereof).
- **Gradients:** Training a global model, where the mini-batch gradients computed by each user are used to define a notion of similarity. This concept is based on the idea that similar data sets will produce similar gradients in expectation over a mini-batch SGD step.

These ideas will be evaluated with respect to the following criteria:

- **User focus:** Freedom of choice left to the users and interpretability of the information presented to them.
- **Resilience:** Compatibility with the distributed machine learning setting, resilience to incomplete data and to unreliably connected users.
- **Incentives:** Potential to provide incentives and/or protect intellectual property.
- **Feedback:** Potential to generate feedback on data quality and feature collection for the user. This does not include feedback on the interoperability of data sets, as this can be done without any data selection method.
- **Privacy:** Compatibility with existing methods that ensure data privacy, such as secure multi-party computation¹ (MPC).
- **Byzantine robustness:** Compatibility with existing methods that ensure robustness against byzantine participants. By their purpose, effective model personalization methods are often robust against the poisoning of other users’ data and highly susceptible to the poisoning of the requesting user’s own data.
- **Complexity:** Computation and communication requirements.

To facilitate their evaluation, ideas will be formalized into concrete methods. The detailed description of these methods, along with their theoretical evaluation with respect to the mentioned criteria, will constitute Section 4 of this report.

¹Put simply, the term *secure multi-party computation* refers to a range of techniques that make it possible for two or more parties to jointly evaluate a function (e.g., `sum` for 3 or more parties) without revealing their individual inputs (cf. Section 4.2.1 of [10]).

2.3 Application case study

Next, we will implement these methods and apply them to a medical data set collected on 577 patients in Sierra Leone during the 2013 - 2015 West African epidemic of Ebola Virus Disease (EVD) [7]. For the sake of reproducibility, we will also use the publicly available Titanic data set [5]. Finally, an application will be made to access the Ebola Data Platform (EDP) assembled by the Infectious Diseases Data Observatory [4].

The cross-silo FL problem will be simulated on a single machine. We will split the available data across simulated agents in a variety of ways to emulate unequal data set size, unequal distribution of features, unequal distribution of labels, unequal distribution of missing entries, and varying degrees of inequality between data sets.

We will select one method among the ones described in Section 4, and evaluate how well it performs for each split by using it to train and evaluate a prediction model. Two baseline methods will also be evaluated, where the model is trained once using only the local data set, and once using all available data. Evaluation metrics, such as root mean squared error or classification accuracy, will be computed on the local test set for the trained models. We will visualize the evolution of these metrics during the learning process, per split for the various methods.

Finally, we will also analyze the parameters of the trained models to understand differences between the methods.

3 Landscape analysis

In this section, we report on tools and literature that are relevant to the project at hand. After discussing recent advancements in privacy-preserving machine learning in Section 3.1, we briefly summarize relevant publications on model personalization in FL in Section 3.2. Subsequently, Section 3.3 introduces the problem of rewarding honest participation, with the notion of *influence* as a potential solution. Finally, available tools to ensure the reproducibility of the work presented in this project are surveyed in Section 3.4.

3.1 Privacy in ML

A very comprehensive overview of the FL problem, including state-of-the-art methods and open research areas of data privacy, resilience, personalization, and robustness, is given by Kairouz et al. [10]. The main tools used to train ML models with data privacy guarantees are MPC², homomorphic encryption³ (HE), and differential privacy⁴ (DP). These tools are implemented in the open-source deep learning framework `PySyft`, which exposes familiar deep learning APIs like `PyTorch` and `TensorFlow` [20, 17].

A different approach is taken in [14]: Here, methods and metrics for privacy-preserving ML and knowledge discovery from data (KDD) are compared and contrasted for each stage of the data lifecycle: data collection, data publication and data output (e.g., in the form of a singular prediction). Strengths and drawbacks of existing DP, MPC, HE, randomization, and k-anonymity tools are discussed from a practical perspective.

As exposed by Kairouz et al., authors often address specific issues individually, such as data privacy *or* robustness [10]. For instance, at the date of publishing, `PySyft` does not provide robustness against malicious users [20]. Nevertheless, a privacy-preserving *and* byzantine-robust⁵ aggregation scheme is proposed in [8]. It relies on securely computing the pairwise distance between the gradients of each client at each round, and using them to select a byzantine-robust subset of clients at each round (e.g. via multi-KRUM [1]).

3.2 Model personalization

Classical ML algorithms often assume that all samples used at training, testing, and prediction time are independent and identically distributed (IID). In the FL setting, however, this assumption can be faulty for various

²For MPC, see footnote in Section 2.2.

³Homomorphic encryption consists of encrypting inputs in such a way that certain mathematical operations can be performed on cyphertext (i.e., on the encrypted version of the inputs) to produce the cyphertext of the correct outputs.

⁴Differential privacy aims to mask the contribution of any individual user to the ML model by adding a small amount of noise, in order to prevent or limit information disclosure. It is discussed in Section 4.2.2 of [10].

⁵Byzantine-robustness is achieved if, given enough participating parties with enough data, the training process converges to a good model even when a given proportion of participants exhibit arbitrarily malicious behaviour (typically up to, but not including, 50%). For instance, instead of following the procedure, these so-called *byzantine participants* can make up arbitrary values whenever they communicate with the server or with other participants.

3.3 Incentives

reasons detailed in Section 3.1 of [10]. One such reason can be that each of the agents participating in the FL task collects data from a different distribution, such as in the example described in Section 2.2 of this report. In such cases, model performance can be improved by applying one or several methods, to which we shall refer as *personalization*. Section 3.3 of [10] outlines the following personalization methods:

Featurization : If samples coming from distinct individuals can be assumed to follow one of several distributions depending on certain personal context information, a natural approach is to include this context information as features in the ML model.

Multi-task learning : Alternatively, each agent’s local problem can be viewed as a separate learning task. If the number of available samples per agent is too low, agents can also be clustered to define one task per cluster.

Local fine-tuning : This method consists of training one global model through FL, broadcasting it, and subsequently letting each agent personalize the model through additional training on their local data set.

Learning-to-Learn (LTL): Also called *Meta-learning*, LTL proposes ways to learn a learning algorithm by sampling from a meta-distribution of ML tasks.

Model-agnostic Meta-learning (MAML): The goal of MAML is to optimize a global model specifically as a good starting point for local fine-tuning. This contrasts the approach of locally fine-tuning a global model that was optimized exclusively for its accuracy, potentially at the expense of its “fine-tuneability”.

As we see, personalization can be achieved by clustering the agents and therefore selecting a subset of the global data set for each task in multi-task learning. Nevertheless, it is not to be confused with the term *training data selection* (TDS), which is generally used to refer to the task of selecting good samples from one large training set whose samples exhibit various levels of noise and relevancy [12].

A learning theoretic formulation of the personalized FL setting, along with three personalization approaches, are presented in [13]. The first presented method, *hypothesis-based user clustering*, is a form of multi-task learning where an appropriate agent clustering is found *during* the model training process. Next up, *data interpolation* aims to reach a good compromise between generalization properties and specificity of the resulting model, by minimizing the loss on a weighted combination of local data and global data. Finally, *model interpolation* is presented as a MAML-like method for the joint optimization of the global model, a global-to-local proportion for each agent, and the local model for each agent. In model interpolation, the purpose of the local models is really to correct wrong predictions of the global model, rather than to make reliable predictions on their own.

3.3 Incentives

Successful cross-silo FL applications depend on the data collection efforts and the subsequent provision of high-quality data by all participating agents. For this reason, designing incentives for honest participation is highlighted as an important practical research question in [10]. Rewarding each agent with money or platform access privileges based on their *influence*, is proposed as a potential solution in [19]. Influence measures how much each agent’s participation has contributed to minimizing the loss function (on a private test set defined by the user), and an efficient approximation scheme is given in [19]. If the user has a sufficiently large private test set, this scheme rewards truthful data collection and reporting, as well as early participation in the training process.

3.4 Reproducible data science

A variety of platforms like Google Colaboratory⁶ and Renku⁷ have been developed to enable reproducible data science by leveraging version control, virtual environments, and containerization, and by providing access to expensive hardware such as GPUs and TPUs⁸. Others, such as Amazon Web Services (AWS)⁹ and Data Version

⁶Google Colaboratory: <https://colab.research.google.com/>

⁷Renku: <https://datascience.ch/renku/>

⁸GPU (graphics processing unit) and TPU (tensor processing unit) are processors which leverage same instruction, multiple data (SIMD) parallelism. The training of most ML models is a heavily data-parallel task, allowing GPUs and particularly the specially-designed TPUs to achieve considerable speed-ups and energy efficiency gains over training on CPUs.

⁹AWS: <https://aws.amazon.com/>

Control (DVC)¹⁰, provide comprehensive services for the entire ML pipeline, from model development and tuning to production. If such services are not used, virtual environments¹¹ are a great tool to ensure the reproducibility of results obtained by executing `Python` programs. Unfortunately, they incur substantial storage space usage due to the duplicate installation of the entire `Python` installation and packages. More parsimonious dependency management tools, such as `watermark`, simply rely on printing software versions in an `IPython` cell [18].

4 Model personalization and data selection

Model personalization can be a solution to the problem of non-IID data in the FL setting. We have seen several personalization approaches in Section 3.2. In this section, we present a formal definition of the task of model personalization, followed by the description and discussion of two novel personalization methods. The cross-silo FL setting and the personalization objective as they are considered in this report are laid out in Section 4.1. Then, Section 4.2 gives a comparative overview over the individual methods introduced in the remainder of this section. Finally, two model personalization methods, the *adapted Ndoye factor* and the gradient-based *Weight Erosion* scheme, are presented and analyzed in Section 4.3 and Section 4.4, respectively.

4.1 Setting and Objective

We consider a network of agents i , each collecting samples $\mathcal{S}_i = \{\mathbf{x}_i^{(n)}, y_i^{(n)}\}_{n=1, \dots, N_i}$ from an underlying distribution \mathcal{D}_i . One agent (agent 0) is called the user and wishes to perform an inference task, such as (regularized) linear or logistic regression, to gain knowledge about the distribution \mathcal{D}_0 from which they collect data. For instance, agent 0 could wish to predict the label y_0^{new} of a new sample after observing its features \mathbf{x}_0^{new} , using some kind of approximation of the conditional probability $p_0(y|\mathbf{x})$.

However, agent 0 knows that the number $|\mathcal{S}_0|$ of samples they have collected is fairly small for the chosen inference task. Luckily, agent 0 believes that *some* of the other agents collect sufficiently interoperable¹² samples \mathcal{S}_i from sufficiently similar underlying distributions \mathcal{D}_i , that the *true loss*¹³ of their model on \mathcal{D}_0 could be reduced by including \mathcal{S}_i in the training process. Unfortunately, agent 0 also knows that this is not the case for *all* agents. Thus, simply including all agents' samples equally in the training process would not help reduce the true loss of agent 0's model.

The task of *model personalization*, in a broad sense, is to give agent 0 a training algorithm which discriminates between the available agents in some way to minimize the true loss of the resulting model on \mathcal{D}_0 .

One class of model personalization methods relies on *data selection*, followed by standard decentralized training on the selected subset of data. The task of data selection is thus to help agent 0 select a subset \mathcal{U}_{train} of agents whose samples will be included in the training process. The optimal subset \mathcal{U}_{train}^* is that which minimizes the true loss of the resulting model on \mathcal{D}_0 .

Model personalization. Formally, we are given:

- A set of agents $i \in \mathcal{U} = \{0, 1, \dots, N\}$
 - Each agent i has collected a set of samples (called *data set*) on the domain $\mathcal{X} \times \mathcal{Y}$: $\mathcal{S}_i = \{\mathbf{x}_i^{(n)}, y_i^{(n)}\}_{n=1, \dots, N_i}$
 - Each agent i collects their samples from an (unknown) underlying distribution \mathcal{D}_i : $(\mathbf{x}_i^{(n)}, y_i^{(n)}) \stackrel{i.i.d.}{\sim} \mathcal{D}_i$

¹⁰DVC: <https://dvc.org/>

¹¹For virtual environments in `Python`, see: <https://docs.python.org/3/tutorial/venv.html>

¹²Within the context of an inference task, two data sets are called *interoperable* if they collect all features chosen for said inference task in the same way. It is sometimes possible to make two data sets interoperable by constructing the required features from other, originally available features. For example, a data set that records the date of symptom onset and the date of admission can be made interoperable with a data set that records the number of days since symptom onset at admission (w.r.t. an inference task for which this number of days is a relevant feature). For a general definition see [3].

¹³The *true loss* of a model on a distribution is formally defined in Equation 1. It is a quantity that can be obtained only analytically, and only if the model and the distribution are perfectly known. In practice, it is approximated by the *test loss* on a test set consisting of samples drawn from the distribution. This approximation is good if the samples of the test set are “sufficiently” numerous and drawn “sufficiently” independently from each other, as well as from the samples of the training set.

- It is assumed that the label y_0 is not independent of the features \mathbf{x}_0 under \mathcal{D}_0 : $p_0(y|\mathbf{x}) \neq p_0(y)$
- A class of models \mathcal{M} s.t. $f : \mathcal{X} \rightarrow \mathcal{Y} \forall f \in \mathcal{M}$. For instance, \mathcal{M} could be the class of linear models where $f(\mathbf{x}) = \mathbf{w}\mathbf{x}$, $\mathbf{w} \in \mathbb{R}^D$.
- A loss function: $\ell(y, \hat{y})$. For instance, the loss function could be the mean squared error (MSE).

We define the true loss $\mathcal{L}_{\mathcal{D}_0}(f)$ of a model $f \in \mathcal{M}$ on \mathcal{D}_0 :

$$\mathcal{L}_{\mathcal{D}_0}(f) = \mathbb{E}_{(\mathbf{x}_0^{new}, y_0^{new}) \sim \mathcal{D}_0} [\ell(y^{new}, f(\mathbf{x}^{new}))] \quad (1)$$

And we attempt to find a training algorithm \mathcal{A} which, given the set \mathcal{U} of agents i with their individual data sets \mathcal{S}_i , the class of models \mathcal{M} , and the loss function ℓ , will produce the model $f = \mathcal{A}(\mathcal{U}, \mathcal{M}, \ell) \in \mathcal{M}$ which minimizes the true loss $\mathcal{L}_{\mathcal{D}_0}(f)$ on \mathcal{D}_0 (Equation 2). This is opposed to standard federated learning methods $\mathcal{A}_{standard}^{FL}$, which consider that all non-malicious agents collect samples from the same underlying distribution. Applying such a method in the personalized FL setting would minimize the loss on a weighted sum of the agents' distributions \mathcal{D}_i , $i \in \mathcal{U}$ instead, as shown in Equation 3. The weights λ_i depend on the specific learning method $\mathcal{A}_{standard}^{FL}$, but they are often either $\lambda_i = \frac{1}{|\mathcal{U}|} \forall i \in \mathcal{U}$ or $\lambda_i = \frac{|\mathcal{S}_i|}{\sum_{i \in \mathcal{U}} |\mathcal{S}_i|}$

$$\textbf{Personalized FL minimization problem:} \quad find : \mathcal{A} = \arg \min_{\mathcal{A}'} \mathcal{L}_{\mathcal{D}_0}(\mathcal{A}'(\mathcal{U}, \mathcal{M}, \ell)) \quad (2)$$

$$\textbf{Standard FL result:} \quad \mathcal{A}_{standard}^{FL} = \arg \min_{\mathcal{A}'} \sum_{i \in \mathcal{U}} \lambda_i \mathcal{L}_{\mathcal{D}_i}(\mathcal{A}'(\mathcal{U}, \mathcal{M}, \ell)) \quad (3)$$

Data selection. The data selection problem is a variant of the model personalization problem, where we are additionally given a training algorithm \mathcal{A} which takes a training set \mathcal{S} of samples as input, and produces a model $f_{\mathcal{S}}$ as output: $f_{\mathcal{S}} = \mathcal{A}(\mathcal{S})$, $\hat{y} = f_{\mathcal{S}}(\mathbf{x})$. For instance, \mathcal{A} could be ridge regression with a regularization parameter selected among a certain number of candidates to minimize the validation error in 4-fold cross-validation.

We introduce the following additional notation:

- A subset \mathcal{U}_{train} of agents: $\mathcal{U}_{train} \subseteq \mathcal{U}$
- The corresponding training set: $\mathcal{S}_{train} = \bigcup_{j \in \mathcal{U}_{train}} \mathcal{S}_j$

And we aim to find the subset \mathcal{U}_{train}^* which will minimize the loss of the trained model $f_{\mathcal{S}_{train}} = \mathcal{A}(\mathcal{S}_{train})$ for agent 0:

$$\mathcal{U}_{train}^* = \arg \min_{\mathcal{U}_{train} \subseteq \mathcal{U}} \mathcal{L}_{\mathcal{D}_0}(f_{\mathcal{S}_{train}})$$

Fictive example. The agents could be individual hospitals dispersed across one or several regions. The aetiology of common, generic symptoms such as fever is highly dependent on geographic location, where the rural setting suffers more faecal-oral and vector borne diseases such as hepatitis A and malaria, while fevers in the urban setting tend to be related to respiratory disease. A medical professional (agent 0) might want to train a ML model to help diagnose the patients admitted to their urban hospital. The number N_0 of samples collected at their hospital, however, is too limited to yield a satisfactory model, so agent 0 considers using data from other hospitals. Knowing that rural populations suffer from fairly different problems than agent 0's urban patients, agent 0 pre-selects only other urban hospitals – in other words, agent 0 performs manual data selection based on prior medical knowledge. However, agent 0 suspects the presence of other confounding variables that could cause the samples collected by *some* of the other urban hospitals to negatively affect the true loss of the trained model on \mathcal{D}_0 . In this example, the underlying distribution \mathcal{D}_0 describes the population of all possible patients of agent 0.

4.2 Summarizing evaluation of methods

In the following discussion of proposed methods, we shall assume without loss of generality that the data sets \mathcal{S}_i collected by each agent i are perfectly interoperable. Indeed, a non-interoperable data set \mathcal{S}_i^- can always

	Ndoeye	Weight Erosion
User focus	o	-
Resilience	+	+
Incentives	+	o
Feedback	+	-
Privacy	+	+
Byzantine	-	+
Complexity	+	o

Table 1: Summarizing evaluation of the model personalization methods outlined in the remainder of this section, along the metrics introduced in Section 2.2. Explanation of the symbols:

- +: This metric is perceived as a (potential) strength of this method.
- o : This metric is viewed as neither a strength, nor a true weakness, of this method.
- : This metric is perceived as a weakness of this method.

be made technically interoperable by filling in all missing features with arbitrary values¹⁴. The resulting data set \mathcal{S}_i^+ no longer truly corresponds to the distribution \mathcal{D}_i , and is likely less useful¹⁵ than it would be if it had directly been sampled from \mathcal{D}_i in an interoperable manner. This loss of usefulness is a consequence of the lack of interoperability in the original data set \mathcal{S}_i^- and can only partially be alleviated by filling in the missing features in an informed way (rather than arbitrarily).

In Section 4.3 and Section 4.4, two novel model personalization methods are presented, and evaluated along the metrics introduced in Section 2.2. Table 1 shows these evaluations, reduced to a general appreciation mark.

4.3 The adapted Ndoeye factor

This method addresses the data selection problem, specifically. It differs from the original Ndoeye factor, in that the available agents i are ranked by increasing \mathcal{D}_i -to- \mathcal{D}_0 transfer loss of the corresponding models $f_i = \mathcal{A}(\mathcal{S}_i)$:

$$\mathcal{L}_i^{transfer} = \mathcal{L}_{\mathcal{D}_0}(f_i) - \mathcal{L}_{\mathcal{D}_i}(f_i)$$

In practice, this transfer loss can only be approximated:

$$\hat{\mathcal{L}}_i^{transfer} = \mathcal{L}_{\mathcal{S}_0}(f_i) - \hat{\mathcal{L}}_{\mathcal{D}_i}^{CV}(f_i)$$

Above, $\mathcal{L}_{\mathcal{S}_0}(f_i)$ represents the average error of f_i on the finite data set \mathcal{S}_0 , while $\hat{\mathcal{L}}_{\mathcal{D}_i}^{CV}(f_i)$ corresponds to the approximation obtained by cross-validation on \mathcal{S}_i . Agent 0 then uses this ranking to select a training set \mathcal{S}_{train} , e.g. by selecting \mathcal{U}_{train} as the k agents with the lowest transfer loss or by selecting all agents whose transfer loss is below a given threshold.

It is clear why it is sensible to rank agents by increasing value of $\mathcal{L}_{\mathcal{D}_0}(f_i)$: If it is low, this most likely implies that \mathcal{S}_i is well-suited as training data for the inference task specified by agent 0. Unfortunately, since $f_i = \mathcal{A}(\mathcal{S}_i)$, $\mathcal{L}_{\mathcal{D}_0}(f_i)$ depends not only on the distribution \mathcal{D}_i (as well as \mathcal{D}_0 and \mathcal{A}), but also on the size $|\mathcal{S}_i|$ (i.e., the number of samples contained in \mathcal{S}_i). However, the size $|\mathcal{S}_i|$ no longer matters when the data sets of several users are combined into \mathcal{S}_{train} . Thus, the influence of $|\mathcal{S}_i|$ is removed by considering the transfer loss $\mathcal{L}_i^{transfer}$ instead of just the loss $\mathcal{L}_{\mathcal{D}_0}(f_i)$.

User focus. This method can be explained intuitively as evaluating how successfully an inference made on \mathcal{S}_i can be transferred to \mathcal{S}_0 . Thus, users can understand the ranking process, albeit without gaining any information as to *why* any particular data set was ranked higher than another. Since the ranking happens before any decentralized training, the user is free to investigate this question by means of further data analysis before manually selecting \mathcal{U}_{train} , or to select \mathcal{U}_{train} automatically.

¹⁴In most cases, one can even fill in the missing features with reasonable, non-arbitrary values. This increases the usefulness of the resulting data set \mathcal{S}_i^+ .

¹⁵In the context of the data selection problem, a data set \mathcal{S}_i is more useful than another data set \mathcal{S}_j , if agent i is more likely to be included in \mathcal{U}_{train}^* than agent j .

Resilience. This method judges each data set as a whole, for which incomplete data is not particularly problematic. Further, unreliably connected agents can be dealt with by using any existing resilient decentralized training scheme. The ranking only requires one message from agent 0 to each other agent and back. The user can be given the option to select a subset \mathcal{U}_{train} and start training before all other agents have responded. However, new agents wishing to join during the training process must first train a local model and be ranked in comparison with the other agents.

Incentives. Agents can be rewarded, either proportionally to their rank, or simply based on whether they were included in \mathcal{U}_{train} . Finally, agent 0 can publish their selected \mathcal{U}_{train} , which would enable the other agents to check whether their contribution is sufficiently diluted to protect their intellectual property. Since this happens before the training process begins, agents can withdraw from any proposed inference tasks that would impinge too much on their intellectual property.

Feedback. If the model weights of f_i are revealed to agent 0, then they can be used to spot differences between \mathcal{S}_0 and \mathcal{S}_i , at the expense of training a model $f_0 = \mathcal{A}(\mathcal{S}_0)$. The other agents can receive extensive feedback based on their ranking, whether or not they were included in \mathcal{U}_{train} , and which model weights were particularly different from f_0 .

Privacy. $\mathcal{L}_{\mathcal{S}_0}(f_i)$ can be computed via MPC without revealing \mathcal{S}_0 to agent i nor f_i to agent 0, albeit at a high communication cost.

Byzantine robustness. This method is not robust against byzantine agents. While a byzantine-robust federated or distributed training algorithm can be used, the user selection process relies on each agent to honestly and correctly report their cross-validation loss $\hat{\mathcal{L}}_{\mathcal{D}_i}^{CV}(f_i)$.

Complexity. In addition to the decentralized training process itself, this method merely requires each agent i to train one model f_i (entirely locally) and to send its weights, along with a singular loss value, back to agent 0. If MPC is used to hide the weights of f_i from agent 0, then the communication requirements increase drastically.

4.4 Gradient-based Weight Erosion

At the core of this method is a novel adaptation of federated training algorithms based on stochastic gradient descent (SGD). It is based on the secure (i.e., privacy-preserving) and byzantine-robust federated SGD-based ML framework proposed in [8]. Like said framework, this method relies on round-based gradient descent, and on securely computing distances between gradients through MPC at each round. The novelty lies in how these distances are used:

At each round, each agent i (including agent 0) computes a gradient:

$$\mathbf{g}_i = \nabla_{\mathbf{w}} \mathcal{L}_{\mathcal{S}_i}(f)$$

Above, \mathbf{w} stands for all learned parameters of the current model f . The (normalized) distance between \mathbf{g}_i and \mathbf{g}_0 is used as a metric for how different \mathcal{D}_i is from \mathcal{D}_0 :

$$d_{i,0}^{rel} = \frac{\|\mathbf{g}_i - \mathbf{g}_0\|}{\|\mathbf{g}_0\|} \geq 0$$

Let us analyze what happens in a few examples, if \mathbf{g}_i are indeed the full gradients (i.e., computed on the entire set \mathcal{S}_i), as opposed to stochastic gradients. Let us further assume that each training loss function $\mathcal{L}_{\mathcal{S}_i}(f)$ is convex in the model parameters \mathbf{w} .

Example 1. Suppose the model has just been randomly initialized to f^0 and performs very sub-optimally for both \mathcal{S}_i and \mathcal{S}_0 . Then, the normalized distance between the gradients will be very low. Formally:

$$\begin{aligned} \text{Let} \quad & \min \{ \mathcal{L}_{\mathcal{S}_0}(f^0), \mathcal{L}_{\mathcal{S}_i}(f^0) \} \gg \max \{ \mathcal{L}_{\mathcal{S}_0}(f_{\mathcal{S}_i}), \mathcal{L}_{\mathcal{S}_i}(f_{\mathcal{S}_0}) \}, \quad f_{\mathcal{S}_i} = \arg \min_{f \in \mathcal{M}} \mathcal{L}_{\mathcal{S}_i}(f) \\ \Rightarrow \quad & \| \mathbf{g}_i - \mathbf{g}_0 \| \ll \| \mathbf{g}_0 \| \\ \Rightarrow \quad & d_{i,0}^{rel} \ll 1 \end{aligned}$$

Therefore, both agents i and 0 are fully included at this stage of the training process, because the distance $d_{i,0}^{rel}$ is close to 0 .

However, as the training process progresses and the model gradually performs better, the distance between the gradients steadily increases. We can analyze the edge cases where f is either the global model or the local model:

Example 2. Suppose f is the global model for 2 agents, 0 and i :

$$\begin{aligned} \text{Let} \quad & f = \arg \min_{f' \in \mathcal{M}} \mathcal{L}_{\mathcal{S}_0 \cup \mathcal{S}_i}(f'), \quad \mathcal{L}_{\mathcal{S}_0 \cup \mathcal{S}_i}(f) = \frac{|\mathcal{S}_0| \mathcal{L}_{\mathcal{S}_0}(f) + |\mathcal{S}_i| \mathcal{L}_{\mathcal{S}_i}(f)}{|\mathcal{S}_0| + |\mathcal{S}_i|} \\ \Rightarrow \quad & \mathbf{0} = \nabla_{\mathbf{w}} \mathcal{L}_{\mathcal{S}_0 \cup \mathcal{S}_i}(f) = \frac{|\mathcal{S}_i|}{|\mathcal{S}_0| + |\mathcal{S}_i|} \left(\frac{|\mathcal{S}_0|}{|\mathcal{S}_i|} \mathbf{g}_0 + \mathbf{g}_i \right) \\ \Rightarrow \quad & \mathbf{g}_i = - \frac{|\mathcal{S}_0|}{|\mathcal{S}_i|} \mathbf{g}_0 \\ \Rightarrow \quad & d_{i,0}^{rel} = \frac{\| \mathbf{g}_i - \mathbf{g}_0 \|}{\| \mathbf{g}_0 \|} = \frac{\| \left(-\frac{|\mathcal{S}_0|}{|\mathcal{S}_i|} - 1 \right) \mathbf{g}_0 \|}{\| \mathbf{g}_0 \|} = 1 + \frac{|\mathcal{S}_0|}{|\mathcal{S}_i|} \end{aligned}$$

In this case, whether (and to which degree) agent i should participate in the training depends on the sizes $|\mathcal{S}_0|$ and $|\mathcal{S}_i|$: Indeed, if $|\mathcal{S}_0| \ll |\mathcal{S}_i|$, then $d_{i,0}^{rel} \approx 1$, indicating that it would be useful to incorporate agent i further in the training process. This is sensible, considering that the much larger number of samples in \mathcal{S}_i could help reduce the generalization error substantially. Inversely, if $|\mathcal{S}_0| \gg |\mathcal{S}_i|$, then it is not useful to include \mathcal{S}_i in training and we are better off only using the (much more numerous) samples collected by agent 0 . Correspondingly, this leads to $d_{i,0}^{rel} \gg 1$.

Example 3. Now let us extend the previous edge case to more agents:

$$\begin{aligned} \text{Let} \quad & f = \arg \min_{f' \in \mathcal{M}} \mathcal{L}_{global}(f'), \quad \mathcal{L}_{global}(f) = \frac{\sum_{i \in \mathcal{U}} |\mathcal{S}_i| \cdot \mathcal{L}_{\mathcal{S}_i}(f)}{\sum_{i \in \mathcal{U}} |\mathcal{S}_i|} \\ \Rightarrow \quad & \mathbf{0} = \nabla_{\mathbf{w}} \mathcal{L}_{global}(f) \\ \Rightarrow \quad & \text{if } \mathbf{g}_0 \neq \mathbf{0}, \text{ then } \exists i \in \mathcal{U} : \langle \mathbf{g}_i, \mathbf{g}_0 \rangle < 0 \\ \Rightarrow \quad & \text{if } \mathbf{g}_0 \neq \mathbf{0}, \text{ then } \exists i \in \mathcal{U} : \| \mathbf{g}_i - \mathbf{g}_0 \| > \| \mathbf{g}_0 \| \\ \Rightarrow \quad & \text{if } \mathbf{g}_0 \neq \mathbf{0}, \text{ then } \exists i \in \mathcal{U} : d_{i,0}^{rel} > 1 \end{aligned}$$

In this setting, 1 is indeed the best lower bound we can prove. To illustrate this, imagine a task with N agents whose gradients are $\mathbf{g}_{i,i>0} = -\frac{1}{N-1} \mathbf{g}_0$ at the global model f . Then $N \rightarrow \infty \Rightarrow \mathbf{g}_i \rightarrow \mathbf{0} \Rightarrow d_{i,0}^{rel} \rightarrow 1$. In that case, the low distance values suggest we should include them all in the training (in other words: stick close to the global model) because this greatly increases the number of available samples.

In general, in this example, relying on the distance measures implies moving towards the model best suited for a weighted subset of data sets (those with $d_{i,0}^{rel} \approx 1$) and neglecting data sets that are in stark disagreement with \mathcal{S}_0 (those with $d_{i,0}^{rel} \gg 1$).

Example 4. Suppose now that f is the local model of agent 0:

$$\begin{aligned}
\text{Let} \quad & f = \arg \min_{f' \in \mathcal{M}} \mathcal{L}_{S_0}(f') \\
\Rightarrow \quad & \mathbf{g}_0 = \nabla_{\mathbf{w}} \mathcal{L}_{S_0}(f) = \mathbf{0} \\
\Rightarrow \quad & \forall i \in \mathcal{U} : d_{i,0}^{rel} \text{ undefined } (+\infty)
\end{aligned}$$

Under the stated assumptions of convexity and full gradient descent, all distances $d_{i,0}^{rel}$ grow without an upper bound as the trained model approaches the local model. In other words, when we move too close to the local model, all data sets start to appear very different from S_0 . Consequently, using any decreasing function of $d_{i,0}^{rel}$ as a similarity metric results in the danger of converging to the local model. One possible strategy to prevent this is to stop training before convergence. Intuitively, training should not be stopped as long as the gradients fit into a D -dimensional cone¹⁶.

Weight Erosion. We propose an adapted SGD scheme in Algorithm 1, which accounts for the concerns highlighted in these examples, and which can be implemented securely using the two-server framework presented in [8]. Unlike the adapted Ndoeye factor, the result of the Weight Erosion scheme depends on the random model initialization in round 0. However, its round-based approach gives it more resilience to randomly appearing or disappearing agents and its use of distances gives it byzantine robustness.

Each agent i is initially given a weight $\alpha_i^0 = 1$. At each round, they compute a mini-batch gradient¹⁷ and their weight α_i is decreased by a small amount $\Delta\alpha_i$ that depends on $d_{i,0}^{rel}$. Thus, all agents' weights decrease over time until they either reach 0 or the maximum number of rounds is reached. Nevertheless, the speed at which α_i decreases depends on how similar S_i is to S_0 .

The batch size is fixed across all agents, such that all samples are treated equally in the first epoch. However, as the number of rounds increases, samples from smaller data sets will be seen more often than samples from larger data sets, because each agent uses the same number of samples per round. This leads to an over-representation of samples from smaller data sets. To counteract this, $\Delta\alpha$ is made to depend also on the average number of times each sample in S_i has been used. The same issue could alternatively be compensated by assigning a second, fixed weight β_i to each agent, which would be unrelated to α_i .

Algorithm 1 relies on four hyper-parameters: The number of rounds r_{max} , the batch size b , the distance penalty factor p_d , and the size penalty factor p_s . The former two are already inherent to SGD, while the penalty factors are introduced to calculate $\Delta\alpha$. One natural choice for the newly introduced parameters would be $p_s = 0$ and $p_d = 0.05$. With these parameters, for example, a gradient \mathbf{g}_i of the same length as \mathbf{g}_0 would lead to a decrease by $\Delta\alpha = \frac{1}{24}$ if it is at a 45° angle with \mathbf{g}_0 , or by $\Delta\alpha = 0.1$ if the gradients are at a 180° angle. The size penalty factor p_s would have to be chosen differently if some agents have much larger data sets than others.

User focus. With this approach, the user does not learn why a certain subset of the available data sets were selected, even if they understand the selection process. They learn nothing beyond "These are the data sets for which $p_i(y|\mathbf{x})$ is similar to the local set". To make things worse, the selection is only revealed to them after the model has been trained.

Resilience. Like the algorithm proposed in [8], Algorithm 1 does not break if users appear or disappear between subsequent rounds. If this is expected, $\Delta\alpha$ can be modified by replacing r with the number of rounds in which agent i has actually participated. When a new agent appears after the first round, their weight should not be initialized to 1, but rather to the mean or median weight of all agents.

Incentives. Agents could be rewarded according to how many rounds they participate in, or according to their weight (summed over all rounds). The examples above give us an intuition that the result of Algorithm 1 lies somewhere between the global model and agent 0's local model. As long as the number and data set size of other users in \mathcal{U} is sufficient to dilute each individual agent's contribution, the algorithm therefore does not

¹⁶If all gradients fit into a cone, then it is possible to improve the loss on all data sets by taking a step in the opposite direction of the axis of the cone.

¹⁷Alternatively, the agents can compute local SGD updates, which are used in the same way.

Algorithm 1: WEIGHT EROSION

Data: A set of agents $i \in \mathcal{U}$, with associated data sets \mathcal{S}_i , a model class \mathcal{M} , and a gradient-based machine learning algorithm \mathcal{A} .

Result: A personalized machine learning model f .

Set a number of rounds r_{max} , a batch size b , a distance penalty factor p_d , and a size penalty factor p_s ;

Initialize $\alpha_i^0 \leftarrow 1 \ \forall i \in \mathcal{U}$;

Randomly initialize the model $f \in \mathcal{M}$;

```

for  $r$  from 1 to  $r_{max}$  do
  for  $i \in \mathcal{U}$ , starting with  $i = 0$  do
    Select a batch of size  $b$  from  $\mathcal{S}_i$  and compute a gradient  $\mathbf{g}_i$  ;
    Compute the distance  $d_{i,0}^{rel}$  ;
     $\Delta\alpha \leftarrow \left(1 + p_s \left\lfloor \frac{(r-1)b}{|\mathcal{S}_i|} \right\rfloor\right) p_d d_{i,0}^{rel}$  ;
     $\alpha_i^r \leftarrow \max\{0, \alpha_i^{r-1} - \Delta\alpha\}$  ;
  end
   $\bar{\mathbf{g}} \leftarrow \frac{\sum_{i \in \mathcal{U}} \alpha_i^r \mathbf{g}_i}{\sum_{i \in \mathcal{U}} \alpha_i^r}$  ;
  Update the model  $f$  based on  $\bar{\mathbf{g}}$  ;
end

```

impinge anyone’s intellectual property. This condition on \mathcal{U} can be verified by each agent before engaging in a round of training.

Feedback. Unfortunately, Algorithm 1 does not produce any information that lends itself to giving feedback on the user’s data collection.

Privacy. A secure (i.e., privacy-preserving) and byzantine-robust SGD protocol based on a framework with two non-colluding servers is proposed in [8], which relies on secret sharing to compute and reveal the distances between gradients through MPC. Weight erosion can be seamlessly integrated into that protocol as an aggregation rule, instead of the byzantine-robust aggregation rule KRUM [1].

Byzantine robustness. Since weight erosion would *replace* the byzantine-robust aggregation rule in the setting of [8], the byzantine robustness property is lost. While byzantine inputs can negatively affect the training process, a byzantine agent i would likely see their weight α_i decline fast.

Complexity. This method relies on computing distances between gradients at each round. However, at each round, only $\mathcal{O}(|\mathcal{U}|)$ distances $\|\mathbf{g}_i - \mathbf{g}_0\|$ are computed, while the byzantine-robust aggregation scheme in [8] computes $\mathcal{O}(|\mathcal{U}|^2)$ pairwise distances $\|\mathbf{g}_i - \mathbf{g}_j\|$ to select a byzantine-robust subset.

5 Application case study

While a purely theoretical examination is sufficient to assess the adapted Ndoeye factor (Section 4.3), the relatively more complicated Weight Erosion scheme from Section 4.4 warrants a practical example to complement its analysis. It is implemented, as briefly discussed in Section 5.1, and applied to two case studies, which are introduced in Section 5.2. Results are finally shown and discussed in Section 5.3.

5.1 Implementation

The Weight Erosion scheme (Algorithm 1) is implemented on top of an early version of the JAX and Haiku based framework for FL simulations built by S. P. Karimireddy [11, 2, 9]. Changes were made to implement the Weight Erosion aggregation rule, to enable loading data sets from local files, and to report performance metrics of the trained model after each communication round. The just-in-time compilation API of JAX is used to speed up the computation of the distances $d_{i,0}^{rel}$.

Third-party reproducible data science platforms, such as those introduced in Section 3.4, are forgone to avoid copying the sensitive medical data sets introduced in Section 2.3 to third-party servers. Instead, all code for this project will be hosted on `GitHub` (available from within the `epfl-iglobalhealth` organization) and executed on the hardware provided by the Machine Learning and Optimization Laboratory (MLO). The repository is cited as a source [6]. The use of virtual environments is forgone in favor of `watermark`, due to their high storage space usage discussed in Section 3.4 [18].

5.2 Case study

Predicting Ebola infection. While our request to access the Ebola Data Platform (cf. Section 2.3) has been granted, the data transfer could not be arranged on time for use within this semester project.

Thankfully, M. Hartley granted access to the raw set of medical data used in [7]. The data set contains 66 features collected on 577 patients, who were admitted at the GOAL-Mathaska Ebola Treatment Center in Port Loko, Sierra Leone, in 2014 and 2015. This data set is used in a FL simulation, in an attempt to learn a prediction model for the Ebola diagnostic (EVD(+)) vs EVD(-)) based on the features used in the triage score proposed by Hartley et al. [7]. These are: *Ebola contact history*, *days since first symptoms*, *conjunctivitis*, *diarrhoea*, *dysphagia*, *haemorrhage*, *fever >38°C*, and *myalgia*. The raw data set is pre-processed following the procedure outlined in [7], mainly to account for missing values. The exact pre-processing steps are documented on `GitHub` [6].

Predicting Survival with the Titanic data set. To make our work reproducible, we also set up a FL simulation using the publicly available Titanic data set [5]. This set collects 14 features on passengers of the cruise ship Titanic, including their survival status. We train a prediction model for the survival status of passengers based on the following features: The *fare* paid by each passenger, their *passenger class*, their *port of embarkation*, whether they travelled *alone or accompanied*, their *sex*, their *age*, and whether they were *adult or minor* (aged 16 or less). The pre-processing procedure is largely aligned with [15], and the exact steps are documented on `GitHub` [6].

Splitting the data set across users. Due to the small size of the Ebola data set compared to typical ML data sets, the number of agents is kept minimal (3 agents, and 4 for the Titanic data set). The samples are apportioned to the agents in the following ways:

AGE_STRICT: Samples are strictly segregated into groups based on their age. For instance, agent 0: patients aged 0 - 20 years old / agent 1: 21 - 40 years old / agent 2: 41+ years old. This split should not introduce any label skew, since the age distribution is the same in the EVD(+) population as in the EVD(-) population, respectively in the survivors and victims of the Titanic accident. We expect a noticeable feature skew, however, since the prevalence of Malaria is strongly correlated with age. In the Titanic data set, the *age*, *adult or minor* and some other features are strongly correlated with age. This feature skew should be sufficient to affect the conditional probability $p(y|\mathbf{x})$.

AGE_SOME: A subset of agents randomly share the totality of a given age group, while other agents contain other age groups exclusively. For instance, agents 0 and 1 randomly partition the patients aged 0 - 40 years old among themselves, while agent 2 has all patients aged 41+ years old.

5.3 Results

The federated training of a prediction model by three agents is simulated on a single machine for the data sets and splits detailed in Section 5.2. Each simulation is repeated three times, such that each agent serves as user once. Every time, the user’s data set is split into a test set and a training set of equal sizes, whereas 100% of the other agents’ data sets are used as training sets. A classification model, consisting of a 2-output linear regression layer followed by log-softmax, is then trained with three different aggregation schemes:

- Global:** The model is trained on all agents’ training sets without Weight Erosion.
- Weight Erosion:** The model is trained on all agents’ training sets with Weight Erosion.
- Local:** The model is trained only on the user’s training set.

5.3 Results

At each communication round, each trained model’s accuracy is measured on the user’s test set and reported along with the weight α_i of each agent in the Weight Erosion scheme.

Predicting Ebola infection with AGE_STRICT split. The lower of the two red dotted lines in each plot of Figure 1 corresponds to the prevalence of EVD(+) in the test set of the user. Contrary to our expectations, it reveals a considerable feature skew, at $p_1(y) < 0.25$ versus $p_0(y) \approx p_2(y) > 0.35$. This might help explain why the global model performs much less well than the local model. In this particular example, the **Weight Erosion** scheme seems to outperform the **Local** scheme, albeit not always (not in Figure 1, right).

In Figure 3, we investigate the weights of the models trained when agent 1 is the user, to understand why the model trained using **Weight Erosion** has a far better test accuracy as seen in Figure 1 (center). Strikingly, the models trained with **Local** and **Weight Erosion** have very similar parameters, while their test accuracies are very different. The **Global** model, on the contrary, has the same classification accuracy on agent 1’s test set as the **Local** model, despite having only negligible weights for all features. Perhaps the more significant bias term of the model trained with **Weight Erosion** is responsible for the improved classification accuracy.

When agent 2 is the user, on the contrary, **Weight Erosion** and **Local** aggregation schemes yield models with identical (and high) test accuracy, whereas only the **Global** aggregation scheme leads to a much less accurate model (cf. Figure 1, right). Figure 4 shows the corresponding model weights. Indeed, in this case, the bias term is nearly identical across **Weight Erosion** and **Local**, and very different for the **Global** aggregation scheme, which supports our hypothesis that the bias term is responsible for the differences observed when agent 1 is the user.

The most important predictors for both agents 1 and 2, seem to be **Quarantine** and **Conjunctivitis**. It is surprising then, that the model trained with the **Global** aggregation scheme gives them small weights *of the opposed sign* in both cases. Figure 2 confirms this same unexpected phenomenon for agent 0.

Predicting Ebola infection with AGE_SOME split. How fast the weight α_i of each agent decreases depends on a number of influences:

Firstly, it depends on how similar the data sets \mathcal{S}_i and \mathcal{S}_0 are with respect to the inference task at hand. This influence causes the weight of each agent to decrease at a different rate. In particular, how similar the data sets are, depends on the underlying distributions \mathcal{D}_i and \mathcal{D}_0 , as well as random chance if $|\mathcal{S}_i|$ and/or $|\mathcal{S}_0|$ are small. Secondly, as discussed in Section 4.4, $d_{i,0}^{rel}$ also depends on how well or how poorly the trained model is currently performing.

And finally, the weights of all agents decrease more rapidly as the number of communication rounds increases, due to the influence of the size penalty factor p_s . The weight of each agent i decreases faster if $|\mathcal{S}_i|$ is smaller. In Figure 5, \mathcal{S}_1 and \mathcal{S}_0 draw their samples from the same population (the age group 0 - 40 years old), thus we have $\mathcal{D}_0 = \mathcal{D}_1$. Disregarding the effect of random chance in the partitioning process, we would therefore expect the weight of agent 2 to decline noticeably faster than that of agent 1 (respectively agent 0), when agent 0 (respectively agent 1) is the user. We would also expect the weights of agents 0 and 1 to remain fairly similar throughout the simulation when agent 2 is the user. In other words, we would expect the dotted lines to be noticeably further apart (with the green dotted line below the other) in Figure 5 (left) and Figure 5 (center), than they are in Figure 5 (right).

On the contrary, we observe that the dotted lines are approximately the same distance apart in all three graphs of Figure 5, and that the green dotted line is above the orange one in Figure 5 (left). It could simply be that $\mathcal{D}_0 = \mathcal{D}_1 \approx \mathcal{D}_2$, which would explain this observation. This implies that we should reconsider our assumption that splitting the data set by patient age affects the conditional probability $p(y|\mathbf{x})$ rather than only the marginal probability $p(y)$. A different explanation would be that the distance $d_{i,0}^{rel}$ depends much more on how well the trained model is currently performing, than on how useful each agent’s data set is in comparison with the other agents’ sets.

Predicting survival of the Titanic cruise with AGE_STRICT split. In Figure 6, we can see the compensating effect of p_s at work, as agent 1’s weight α_1 experiences the slowest decrease in all three graphs because

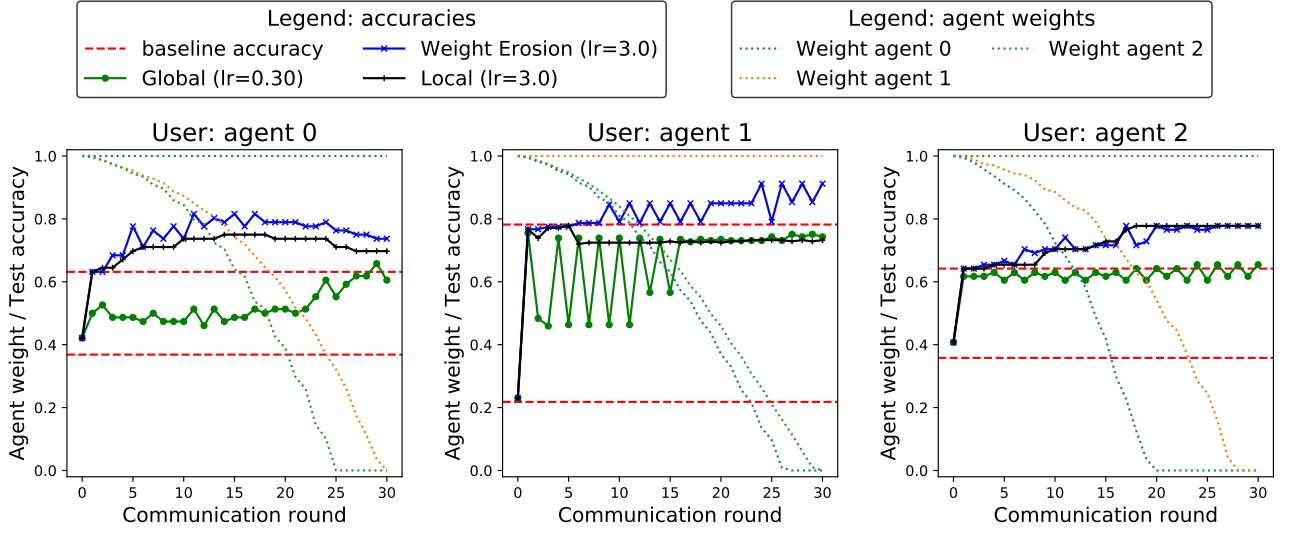


Figure 1: Predicting Ebola infection with AGE_STRICT split. The Ebola data set is split across 3 agents by age group:

agent 0: 0 - 20 years old / agent 1: 21 - 40 years old / agent 2: 41+ years old.

Full lines represent each model's accuracy on the user's test set. In red (dashed), the accuracy obtained by predicting always EVD(+) or always EVD(-). The pointed lines represent the weight α_i of each agent in the Weight Erosion scheme. The learning rates were tuned independently for each aggregation scheme (Global, Weight Erosion, Local), and they are displayed in the legend. One batch of 125 samples per round and agent. Each agent's data set contains two such batches.

$p_d = 0.01$, $p_s = 0.2$, seed = 278.

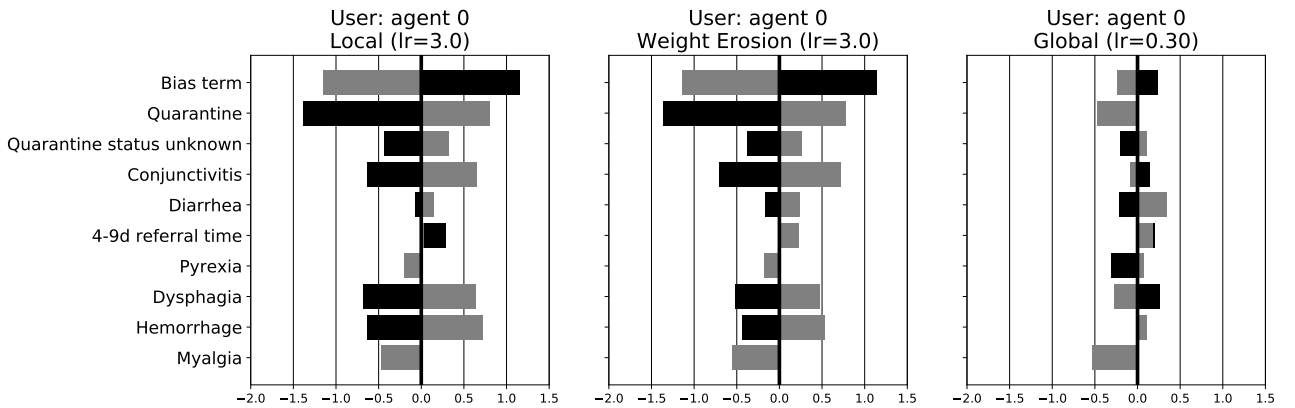


Figure 2: Parameters of the models in Figure 1 (left).

■ : Weight of the feature (or bias) for the first output of the linear regression model.
 ■ : Weight of the feature (or bias) for the second output of the linear regression model.
 <<<<<< >>>>>> : The feature is strongly correlated with EVD(-).
 <<<<<< >>>>>> : The feature is strongly correlated with EVD(+).

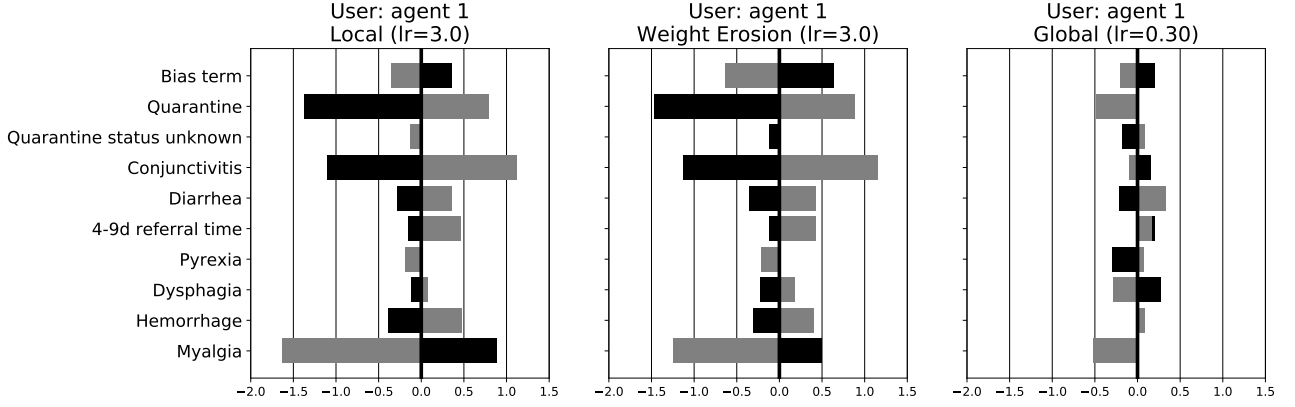


Figure 3: Parameters of the models in Figure 1 (center).

■ : Weight of the feature (or bias) for the first output of the linear regression model.
 ■ : Weight of the feature (or bias) for the second output of the linear regression model.
 <<<<<< >>>>>> : The feature is strongly correlated with EVD(-).
 <<<<<< >>>>>> : The feature is strongly correlated with EVD(+).

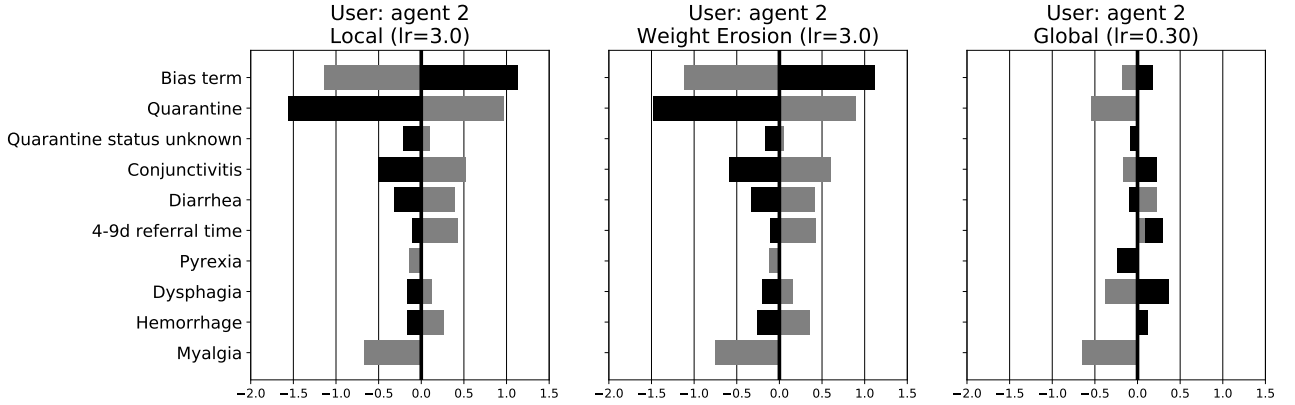


Figure 4: Parameters of the models in Figure 1 (right).

■ : Weight of the feature (or bias) for the first output of the linear regression model.
 ■ : Weight of the feature (or bias) for the second output of the linear regression model.
 <<<<<< >>>>>> : The feature is strongly correlated with EVD(-).
 <<<<<< >>>>>> : The feature is strongly correlated with EVD(+).

5.3 Results

agent 1 has collected one batch more than the other agents. Nevertheless, the difference between α_1 and the other agent weights is only very marked when agent 2 is the user, showing that the impact of p_d rivals that of p_s .

When splitting the Titanic data set into clients by age, we notice that the survival rate (Figure 6, lower of the two red dotted lines in each graph) is fairly similar across all age groups. This implies that the *age* feature may not be particularly useful in predicting the survival of passengers. However, the survival rate is much lower among the group of passengers whose age is unknown. This could simply be a sampling bias, as it would seem easier to research the age of survivors than that of the deceased.

Another debatable feature is whether the passenger is adult or minor (**Passenger is minor**), which is constant across the data sets of agents 1 and 2, and unknown for all passengers in the set of agent 3. We observe that agent 0 stands out by the fact that the **Weight Erosion** model clearly out-performs the **Local** model. Nevertheless, the difference in prediction accuracy does not come from the **Local** model overfitting on the **Passenger is minor** feature, since its weights are nearly identical in both models as we see in Figure 7.

Further, we notice that **Travelling alone** is strongly correlated with death in both models, indicating that passengers aged 0 - 20 years old were more likely to die if they were travelling alone. Notably, the same feature is weakly correlated with survival (rather than death) in the **Global** model, as well as in the very accurate **Local** and **Weight Erosion** models trained with agent 2 as user (cf. Figure 8, left and center).

Three features have smaller weights in the **Weight Erosion** model than in the **Local** model: **Travelling alone**, boarding in **Queenstown** (as opposed to Southampton or Cherbourg), and travelling in **Second Class**. These differences could be a sign that the local model is overfitting, especially if they apply to a small portion of the passengers, as is the case with **Queenstown** and presumably **Travelling alone**.

Finally, we are very interested in noting that the **Weight Erosion** model is far from being a weighted average of the **Global** and **Local** models. Quite on the contrary, the weights of features such as **Sex**, **Age**, or boarding in **Queenstown**, differ more between the **Global** model and the **Weight Erosion** model, than they do between the **Global** model and the **Local** model.

In Figure 8 (left and center), we investigate the weights of the models trained with **Local** and **Weight Erosion** schemes when agent 2 is the user, since these models perform exceptionally well ($> 90\%$) despite the unexceptional survival probability of 40%. We observe that weights with absolute values < 0.25 differ between the models, without affecting their performance, while larger weights are very consistent.

Further, two main features stand out: Firstly, the weights for **Sex** are spectacularly large (with absolute values between 1.1 and 1.4). Given the models' extraordinary test accuracies, we conclude that most survivors aged 36+ were ladies. Secondly, the weights of the **First Class** and **Second Class** features show that, in this age group, only First-Class passengers were much more likely to survive than Third-Class passengers.

Predicting survival of the Titanic cruise with AGE.SOME split. Unlike with the AGE.SOME split of the Ebola data set, Figure 9 matches our expectations. Indeed, we observe that agents 0 and 1, whose samples are drawn from the same age group, have the highest weight in each other's model, while their weights are similar to each other in the two other agents' models.

The models trained by these agents achieve the same test accuracy as with the AGE.STRICT split, save for the spectacular accuracy obtained for agent 2 in Figure 6.

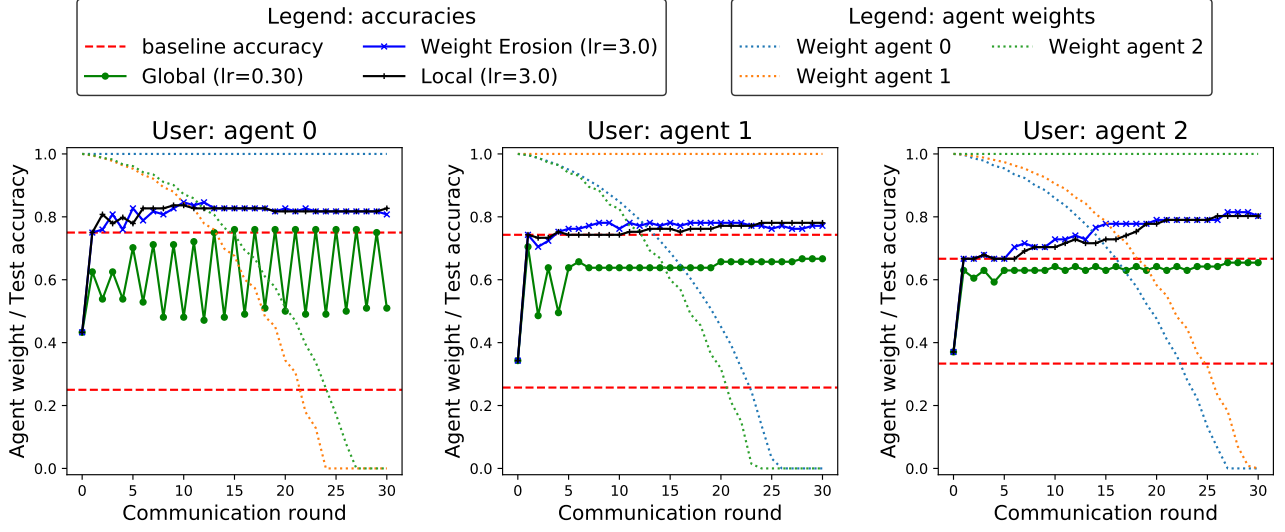


Figure 5: Predicting Ebola infection with AGE.SOME split. The Ebola data set is split across 3 agents by age group:

agents 0 and 1: 0 - 40 years old / agent 2: 41+ years old.

Full lines represent each model's accuracy on the user's test set. In red (dashed), the accuracy obtained by predicting always EVD(+) or always EVD(-). The pointed lines represent the weight α_i of each agent in the Weight Erosion scheme. The learning rates are displayed in the legend. One batch of 125 samples per round and agent. Each agent's data set contains two such batches.

$p_d = 0.01$, $p_s = 0.2$, seed = 278.

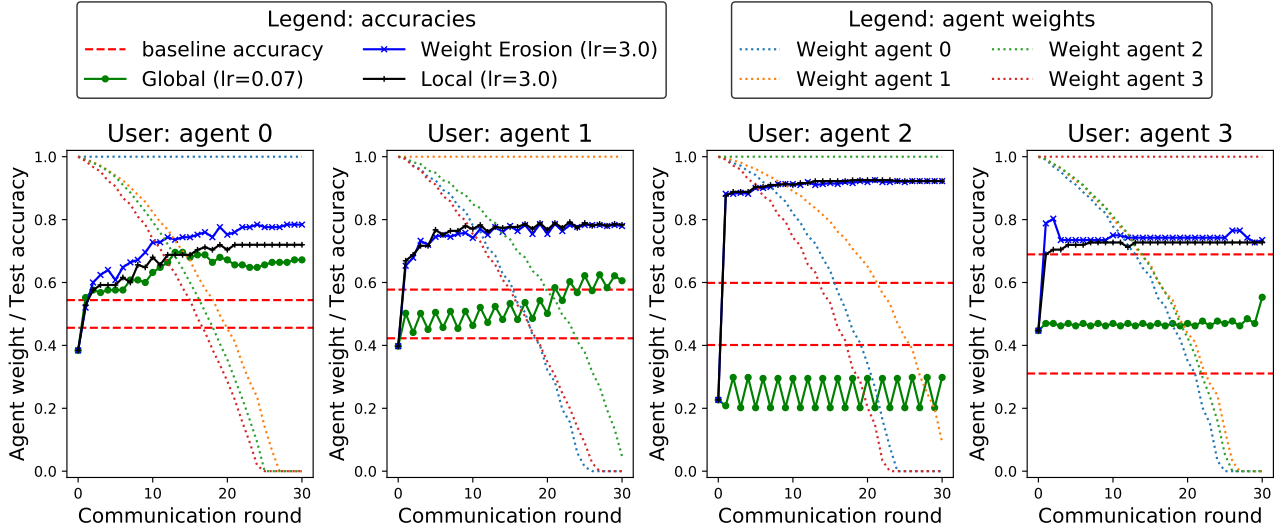


Figure 6: Predicting survival of the Titanic cruise with AGE.STRICT split. The Titanic data set is split across 4 agents by age group:

agent 0: 0 - 20 years old / agent 1: 21 - 35 years old / agent 2: 36+ years old / agent 3: age unknown.

Full lines represent each model's accuracy on the user's test set. In red (dashed), the accuracy obtained by predicting always 1 or always 0. The learning rates were tuned independently for each aggregation scheme (Global, Weight Erosion, Local), and they are displayed in the legend. One batch of 161 samples per round and agent. Agent 1 has 3 such batches, other agents have 2 each.

$p_d = 0.01$, $p_s = 0.2$, seed = 278.

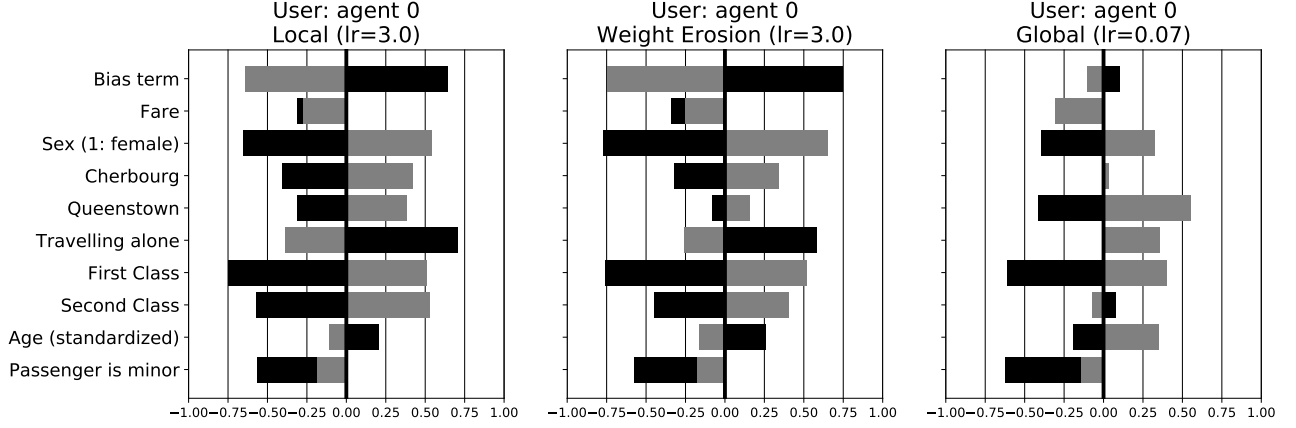
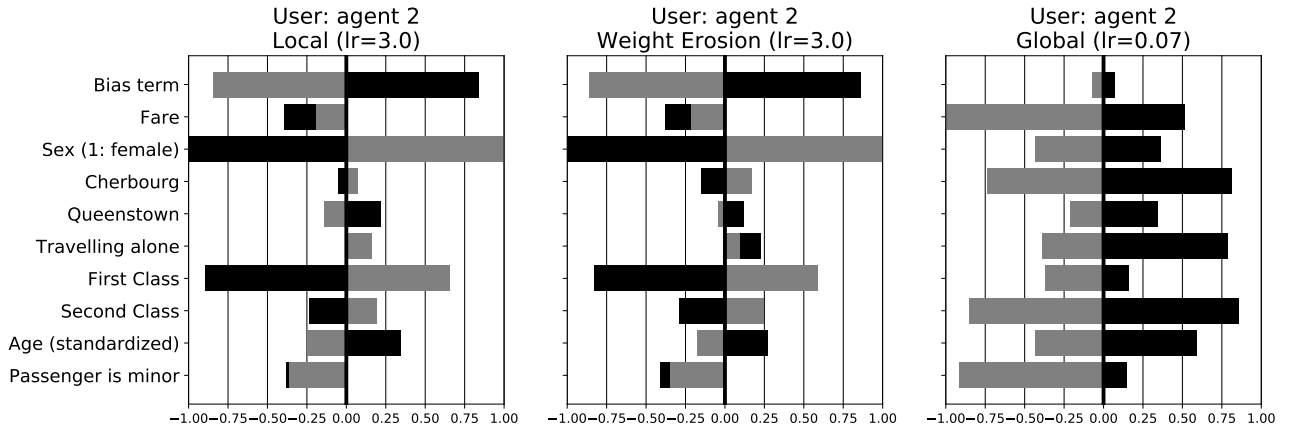


Figure 7: Parameters of the models trained in Figure 6 (left).

: Weight of the feature (or bias) for the first output of the linear regression model.
 : Weight of the feature (or bias) for the second output of the linear regression model.
 <<<<<< >>>>>> : The feature is strongly correlated with death.
 <<<<<< >>>>>> : The feature is strongly correlated with survival.

Figure 8: Parameters of the models in Figure 6 (center right). The test accuracy of the Local and Weight Erosion models is particularly good ($< 90\%$).

: Weight of the feature for the first output of the linear regression model.
 : Weight of the feature for the second output of the linear regression model.
 <<<<<< >>>>>> : The feature is strongly correlated with death.
 <<<<<< >>>>>> : The feature is strongly correlated with survival.

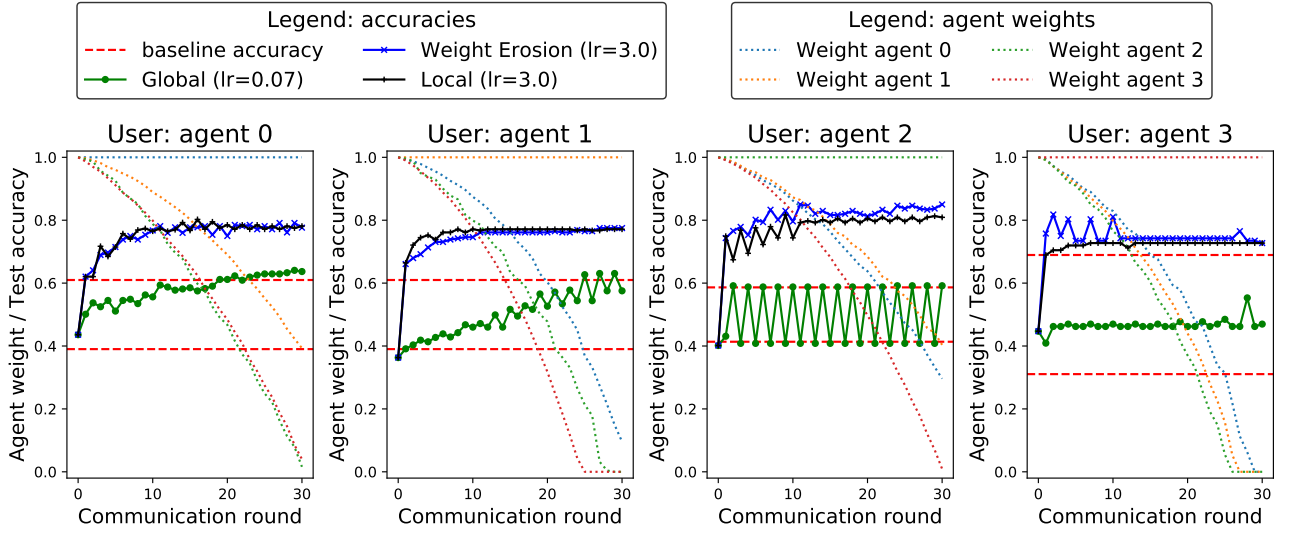


Figure 9: Predicting survival of the Titanic cruise with `AGE.SOME` split. The Titanic data set is split across 4 agents by age group:

agents 0 and 1: 0 - 35 years old / agent 2: 36+ years old / agent 3: age unknown.

Full lines represent each model's accuracy on the user's test set. In red (dashed), the accuracy obtained by predicting always 1 or always 0. The learning rates were tuned independently for each aggregation scheme (Global, Weight Erosion, Local), and they are displayed in the legend. One batch of 132 samples per round and agent. Agent 3 has 2 such batches, other agents have 3 each.

$p_d = 0.01$, $p_s = 0.2$, seed = 278.

6 Conclusion

We introduced, analyzed, and compared, two novel model personalization methods in the cross-silo FL setting: the adapted Ndoye factor, and the Weight Erosion aggregation scheme. We implemented the latter and compared it to two baseline aggregation schemes in two case studies: Firstly, replicating a diagnostic model from a medical science publication, and, secondly, training a classification model on a publicly available data set. We demonstrated that the novel Weight Erosion scheme can outperform both baseline aggregation schemes in some cases. While our analysis exposed the danger of converging to the local model, the case studies revealed that better points of convergence exist and are often reached.

Further theoretical work. Additional refinement is needed to address the under-representation of samples from larger data sets in a more equitable way. Indeed, with the Weight Erosion scheme proposed here, samples from small data sets are still over-represented in the first training rounds. While this is compensated by allowing larger data sets to remain relevant further into the training process, it would be more equitable to solve the issue with a different weight β_i that is unrelated to α_i .

Equally importantly, the influence-based model personalization idea sketched in Section 2.2 should be developed into a formal method. It should further be compared with the adapted Ndoye factor, the Weight Erosion aggregation scheme, as well as methods from the literature, along the criteria presented in Table 1. Finally, the notion of *similarity*, should be defined in mathematical terms, possibly using the notion of *discrepancy* between distributions [13].

Further practical work. The presented application case study would benefit from automated hyper-parameter tuning via cross-validation, as well as from splitting the data sets across varying numbers of users in more diverse ways to generate intentional label skew, $p(\mathbf{x}|y)$ skew, $p(y|\mathbf{x})$ skew, or IID data sets with no systematic skew. The case study should also be extended to other tasks, such as prediction of continuous variables, or image recognition with neural networks. Finally, the examined methods should also be compared in a range of standard tests like training a model within a limited number of SGD steps or a limited communication budget.

References

- [1] Peva Blanchard, Rachid Guerraoui, Julien Stainer, et al. “Machine learning with adversaries: Byzantine tolerant gradient descent”. In: *Advances in Neural Information Processing Systems*. 2017, pp. 119–129.
- [2] James Bradbury et al. *JAX: composable transformations of Python+NumPy programs*. Version 0.1.55. 2018. URL: <http://github.com/google/jax>.
- [3] *Definition: Interoperability*. URL: <http://interoperability-definition.info/en/> (visited on 06/01/2020).
- [4] *Ebola Data Platform — Infectious Diseases Data Observatory*. URL: <https://www.iddo.org/ebola/data-sharing/accessing-data> (visited on 06/11/2020).
- [5] Thomas Cason Frank E. Harrell Jr. *Titanic dataset*. Oct. 2017. URL: <https://www.openml.org/d/40945>.
- [6] Felix Hans Michel Grimberg and Martin Jaggi. *Semester project on private and personalized ML*. 2020. URL: <https://github.com/epfl-iglobalhealth/coML-Personalized-v2>.
- [7] Mary-Anne Hartley et al. “Predicting Ebola infection: A malaria-sensitive triage score for Ebola virus disease”. In: *PLoS neglected tropical diseases* 11.2 (2017).
- [8] Lie He, Sai Praneeth Karimireddy, and Martin Jaggi. “Secure Byzantine-Robust Machine Learning”. In: *arXiv:2006.04747 [cs, stat]* (June 8, 2020). arXiv: 2006.04747. URL: <http://arxiv.org/abs/2006.04747> (visited on 06/16/2020).
- [9] Tom Hennigan et al. *Haiku: Sonnet for JAX*. Version 0.0.1. 2020. URL: <http://github.com/deepmind/dm-haiku>.
- [10] Peter Kairouz et al. “Advances and Open Problems in Federated Learning”. In: *arXiv:1912.04977 [cs, stat]* (Dec. 10, 2019). arXiv: 1912.04977. URL: <http://arxiv.org/abs/1912.04977> (visited on 06/04/2020).
- [11] Sai Praneeth Karimireddy. *JAX federated learning*. 2020. URL: <https://tinyurl.com/Karimireddy-Jax-FL>.
- [12] Miaofeng Liu et al. “Reinforced Training Data Selection for Domain Adaptation”. In: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics. Florence, Italy: Association for Computational Linguistics, 2019, pp. 1957–1968. DOI: 10.18653/v1/P19-1189. URL: <https://www.aclweb.org/anthology/P19-1189> (visited on 03/27/2020).
- [13] Yishay Mansour et al. “Three Approaches for Personalization with Applications to Federated Learning”. In: *arXiv:2002.10619 [cs, stat]* (Feb. 24, 2020). arXiv: 2002.10619. URL: <http://arxiv.org/abs/2002.10619> (visited on 06/03/2020).
- [14] Ricardo Mendes and João P Vilela. “Privacy-preserving data mining: methods, metrics, and applications”. In: *IEEE Access* 5 (2017), pp. 10562–10582.
- [15] Baligh MNassri. *Titanic: logistic regression with python*. 2020. URL: <https://www.kaggle.com/mnassrib/titanic-logistic-regression-with-python?scriptVersionId=26445092>.
- [16] Mohamed Ndoye et al. “Collaborative privacy”. In: (2020). URL: <https://www.mndoye.com/collaborativeprivacy.pdf>.
- [17] OpenMined et al. *PySyft: A library for encrypted, privacy preserving machine learning*. 2020. URL: <https://github.com/OpenMined/PySyft>.
- [18] Sebastian Raschka. *watermark: An IPython magic extension for printing date and time stamps, version numbers, and hardware information*. Version 2.0.2. 2019. URL: <https://pypi.org/project/watermark/>.
- [19] Adam Richardson, Aris Filos-Ratsikas, and Boi Faltings. “Rewarding High-Quality Data via Influence Functions”. In: *arXiv preprint arXiv:1908.11598* (2019).
- [20] Theo Ryffel et al. “A generic framework for privacy preserving deep learning”. In: *arXiv:1811.04017 [cs, stat]* (Nov. 13, 2018). arXiv: 1811.04017. URL: <http://arxiv.org/abs/1811.04017> (visited on 06/13/2020).