

# The Applicability of Federated Learning to Official Statistics

Joshua Stock<sup>1</sup>, Oliver Hauke<sup>2</sup>, Julius Weißmann<sup>2</sup>, and Hannes Federrath<sup>1</sup>

<sup>1</sup> Universität Hamburg, Hamburg, Germany

<sup>2</sup> Federal Statistical Office (Destatis), Wiesbaden, Germany

**Abstract.** This work investigates the potential of Federated Learning (FL) for official statistics and shows how well the performance of FL models can keep up with centralized learning methods. FL is particularly interesting for official statistics because its utilization can safeguard the privacy of data holders, thus facilitating access to a broader range of data. By simulating three different use cases, important insights on the applicability of the technology are gained. The use cases are based on a medical insurance data set, a fine dust pollution data set and a mobile radio coverage data set – all of which are from domains close to official statistics. We provide a detailed analysis of the results, including a comparison of centralized and FL algorithm performances for each simulation. In all three use cases, we were able to train models via FL which reach a performance very close to the centralized model benchmarks. Our key observations and their implications for transferring the simulations into practice are summarized. We arrive at the conclusion that FL has the potential to emerge as a pivotal technology in future use cases of official statistics.

## 1 Introduction

The aim of national statistical offices (NSOs) is to develop, produce and disseminate high-quality official statistics that can be considered a reliable portrayal of reality [Yun+22]. In order to effectively capture our rapidly changing world, NSOs are currently undergoing a process of modernization, leveraging new data sources, methodologies and technologies.

NSOs have effectively extracted information from new data sources, such as scanner data<sup>3</sup> or Mobile Network Operator (MNO) data<sup>4</sup>. However, the potential of numerous other data sources, including privately held data<sup>5</sup> or data from certain official entities, remains largely untapped. Legal frameworks, which are fundamental to official statistics, only adapt slowly to changing data needs and currently hinder access to valuable new data sources. Cooperation with potential data donors faces restrictions due to concerns about privacy, confidentiality, or disclosing individual business interests.

In the meantime, the methodology employed by NSOs is evolving, with machine learning (ML) gaining substantial popularity and, as a result, undergoing a process of establishment. ML has been applied in various areas of official statistics (e.g. [DB17; BDF18; Eur22]), and new frameworks such as [Yun+22] address the need to measure the quality of ML.

Within official statistics, ML tools have proven effective in processing new data sources, such as text and images, or enabling the automation of statistical production tasks, including classifying information or predicting not (yet) available data.

**Federated learning (FL)** is an emerging approach within ML that provides immense unexplored potential for official statistics. It addresses the challenge of extracting and exchanging valuable global information from new data sources without compromising the privacy of individual data owners. Introduced in [McM+17], FL enables collaborative model training across distributed data sources while preserving data privacy by keeping the data localized. In scenarios where external partners are unwilling to share individual-level information due to regulatory or strategic considerations, but still aim to analyze or disseminate global insights in their field of application, NSOs can offer trustworthy solutions by utilizing FL. In return, FL empowers contributing NSOs to integrate new data sources into statistical production.

<sup>3</sup> Scanner data in consumer price statistics and for determining regional price differences <https://www.destatis.de/EN/Service/EXSTAT/Datensaetze/scanner-data.html>, accessed on July 17, 2023

<sup>4</sup> Use of MNO data [https://cros-legacy.ec.europa.eu/content/12-use-mno-data\\_en](https://cros-legacy.ec.europa.eu/content/12-use-mno-data_en), accessed on July 17, 2023

<sup>5</sup> Guidance on private sector data sharing <https://digital-strategy.ec.europa.eu/en/policies/private-sector-data-sharing>, accessed on July 17, 2023

Although FL has been successfully applied to many domains, to the best of our knowledge, besides our work only one currently presented study investigates the applicability of FL to the field of official statistics. In a proof of concept (PoC) by the United Nations (UN), FL is applied to estimate human activity based on data collected from smart and wearable devices [Tem22; Buc23]. The PoC emphasizes operative aspects of FL coordinating multiple NSOs and benefits of additional privacy enhancing technologies.

The main contribution of this paper lies in presenting three additional applications of FL that address current data need representative for official statistics. Complementary, we emphasize measuring the numerical predictive performance and reproducibility by openly sharing our code, which, in two instances, is applied to publicly available data. In the first simulation related to health, individual healthcare costs are predicted utilizing tools for regression. In the second simulation related to **sustainability**, current fine dust pollution levels are classified based on meteorological data. In the third simulation related to **mobility**, the daily range of movement of mobile phone users are classified by MNO data. The first two simulations focus on assessing the estimation performance achieved by FL in comparison to centralized models that have complete access to all available data. The third application presents valuable insights and lessons learned from the implementation of FL, involving the active participation of a real external partner. We draw conclusions on the applicability of FL in NSOs in [section 5](#), which are summarized in [section 6](#).

## 2 Background

Before presenting the simulated use cases in [section 3](#), this section provides an overview of FL and privacy challenges with ML.

### 2.1 Federated Learning

In FL, a centralized server (or aggregator, in our case a NSO) coordinates the process of training a ML model (mainly deep neural networks) by initializing a global model and forwarding it to the data owners (clients). In each training round, each client trains the model with their private data and sends the resulting model back to the central server. The central server uses a FL method to aggregate the updates of the participants into the next iteration of the global model and starts the next round by distributing the updated model to the clients. This process is repeated to improve the performance of the model.

NSOs primarily strive to generate global models that accurately represent the available data, which, in our setting, is distributed among multiple clients. Thus, we compare the performance of FL to models with access to the combined data of all clients. Alternatively, if upcoming applications seek to supply each client with an optimized individual model by leveraging information from the other clients, *personalized* FL can be used. This approach is not covered in this paper but can be found in [KKP20; Hu+18].

### 2.2 Privacy Challenges with Machine Learning

When training data for a ML model is distributed among multiple parties, the data traditionally needs to be combined on a central server prior to training an ML model. FL has become a popular alternative to this approach, as it allows to train a model in a distributed way from the start, without the need to aggregate training data first. Thus, using FL has the privacy advantage that there is no need to exchange private training data. Instead, data holders can train a global model collaboratively in a distributed fashion, without transferring any data record.

But although FL makes sharing private training data obsolete, there are other privacy challenges inherent to ML which have also been observed for FL. While ML models are always trained to fulfill a dedicated task, often more information than strictly necessary for fulfilling the task is extracted into the model weights during training [SRS17]. This excessive, and potentially private, information in the model weights is called privacy leakage. In general, this leakage can be leveraged by any party who has full access to a model and its trained weights.

One concrete example of such a privacy attack is *training data extraction* [ZLH19], which allows extracting data records from a trained model. Another known attack is *model inversion* [HAP17], where repeated requests to the model are used to reconstruct class representatives. *Membership*

*inference* [Sho+17] aims at individual training data records: the attack’s target is to decide whether a specific data record was part of the training data. Building on the original proposal, other works have transferred membership inference attacks to the FL scenario [NSH19]. Last but not least, *property inference* attacks [Mel+19] allow to deduce statistical properties of the target model’s training data. This is especially relevant in FL scenarios, where the characteristics of each client’s local data set can be highly sensitive, e.g., in medical domains.

The applicability of these attacks depends on the concrete use case, the type of model and other factors. Concerning attacker models, i.e., the scenario in which an attack is executed, some FL-specific attacks rely on a malicious aggregator. Nonetheless, all attacks mentioned above also work in an environment where not the aggregator, but one of the FL clients is the attacker. Hence, even if the aggregator can be trusted, e.g., because the aggregator’s role is assumed by a NSO, these attacks can still be executed by other FL clients. Analyzing the individual privacy leakage of the simulated use cases in this paper are out of scope. Nonetheless, raising awareness to these issues, e.g., by communicating potential risks to clients in an FL scenario, should not be neglected. Beyond this, strategies under the umbrella term *privacy-preserving machine learning* (PPML) can help to mitigate these risks [YZH21].

### 2.3 Frameworks

In our simulations, we use the frameworks TensorFlow<sup>6</sup> for neural networks and TensorFlow Federated<sup>7</sup> for FL. We use PyCaret<sup>8</sup> for automizing benchmark experiments in the centralized settings and scikit-learn<sup>9</sup> for data processing.

The code we have written for this work is openly available on GitHub<sup>10</sup>.

## 3 Simulations

Most relevant for NSOs is *cross-silo* FL, where a few reliable clients train a model, e.g. official authorities. In contrast, *cross-device* FL uses numerous clients, e.g. smartphones, to train a model. To analyze the potential of cross-silo FL for official statistics, we run simulations with three different data sets. For each use case, we first compute benchmarks by evaluating centralized ML models, i.e., models which are trained on the whole data set. Afterwards, we split the data set and assign the parts to (simulated) FL clients for the FL simulation. This way, we have a basis for interpreting the performance of the model resulting from the FL training simulation. The performance metrics of the trained ML models (including coefficient of determination  $R^2$  or accuracy) are computed on test sets of each data set.

### 3.1 Medical insurance data

The demand for timely and reliable information on public health is steadily increasing. The COVID-19 pandemic has significantly accelerated this trend, raising questions about the financial feasibility of our healthcare system and the availability of medical supplies.

Thus, our first experiment focuses on modeling a regression problem related to healthcare by considering the following question: Given an individual’s health status characteristics, what is the magnitude of their insurance *charges*? We aim to address two primary questions. Firstly, we explore the suitability of neural networks in comparison to other models for the regression task. Secondly, we assess the feasibility of utilizing a simulated decentralized data set in an FL setting to tackle the problem.

<sup>6</sup> TensorFlow <https://www.tensorflow.org/>, accessed on July 17, 2023

<sup>7</sup> TensorFlow Federated: Machine learning on decentralized data <https://www.tensorflow.org/federated>, accessed on July 17, 2023

<sup>8</sup> PyCaret <https://pycaret.org/>, accessed on July 17, 2023

<sup>9</sup> scikit-learn, Machine Learning in Python <https://scikit-learn.org/>, accessed on July 17, 2023

<sup>10</sup> Code repository for this paper: <https://www.github.com/joshua-stock/fl-official-statistics>, accessed on July 17, 2023. Note that for the mobile radio coverage simulation, the code has only been executed locally on the private data set, hence it is not included in the repository.

*Data set* The given data set links medical insurance premium *charges* to related individual attributes<sup>11</sup>. Considered are the six features *age*, *sex*, *bmi* (body mass index), *children* (count), *smoker* (yes/no) and four *regions*. In our studies, the feature *region* was excluded during FL training and solely utilized for partitioning the data within the FL setting. In total, the data set consists of 1338 complete records, i.e. there are no missing or undefined values. Also, the data set is highly balanced: The values in *age* are evenly dispersed, just as the distribution of male and female records is about 50/50 (attribute *gender*) and each *region* is represented nearly equally often. The origin of the data is unknown, however its homogeneity and integrity suggest that it has been created artificially.

*Data preprocessing* We encode the binary attributes *sex* and *smoker* into a numeric form (0 or 1). The attributes *age*, *bmi* and *children* are scaled to a range from 0 to 1. In the centralized benchmarks, the attribute *region* is one-hot-encoded.

*Setup* We aim to investigate the suitability of neural networks for estimating insurance *charges* and explore the extent to which this problem can be addressed using a FL approach. To achieve this, we compare different models and evaluate their performance.

A basic fully connected neural network architecture, that takes five input features, is utilized. The network consists of three hidden layers with 40, 40, and 20 units in each respective layer. Following each layer, a Rectified Linear Unit (ReLU) activation function is applied. The final output layer comprises a single neuron. To optimize the network, the Adam optimizer with a learning rate of 0.05 is employed. In the federated setting, we utilize the same initial model but integrate FedAdam for server updates. This decision is based on previous research [Red+20], which emphasizes the benefits of adaptive server optimization techniques for achieving improved convergence.

In the centralized approach, we allocate a training budget of 100 epochs. In contrast, the federated approach incorporates 50 rounds of communication between the client and server during training. Each round involves clients individually training the model for 50 epochs. To track the running training, 10% evaluation data is used by each client in the FL setting and 20% is used in the centralized scenario. It is neglected in calculating the final test performance. The remaining shallow learning models undergo hyperparameter optimization using a random search approach with a budget of 100 iterations. We evaluate all models using 5-fold cross validation.

Model	R <sup>2</sup> (± std)	Rel. loss (%)
neural network	81.5(4.01)	3.5
neural network (federated)	78.4(3.13)	7.2
random forest	84.5(4.73)	0.0
XGBoost	84.3 (3.96)	0.2
decision tree	84.1 (4.23)	0.5
k-nearest neighbors	74.4 (5.53)	12.0
linear regression	72.8 (6.07)	13.8

**Table 1.** Performance comparison of different prediction models for the medical insurance use case. The performance is quantified using R<sup>2</sup> in %, along with the corresponding standard deviation (std). Additionally, the relative loss to the best centralized model (rel. loss) is reported.

*Results* We conduct a performance comparison of the models based on their 5-fold cross-validation R<sup>2</sup> scores and consider their standard deviation (see Table 1). The random forest model achieves the highest performance with an R<sup>2</sup> of 84.5 %, closely followed by XGBoost and Decision Tree, which scores 0.2 and 0.5 percentage points lower, respectively.

<sup>11</sup> US health insurance dataset <https://www.kaggle.com/datasets/teertha/ushealthinsurancedatas>, accessed on July 17, 2023

The neural network model achieves an  $R^2$  of 81.5 %, indicating a performance 3.5 % worse than the best model. However, it still provides a reasonable result compared to K-Nearest Neighbors (KNN) and Linear Regression, which obtain significantly lower  $R^2$  scores of 12 % and 13.8 %, respectively.

The Federated neural network demonstrates an  $R^2$  of 78.4 %, slightly lower than the centralized neural network but 7.2 % worse than the random forest model. Notably, the Federated neural network exhibits a lower standard deviation of 3.99 compared to the centralized neural network (4.92) and also outperforms the random forest model (4.73) in this regard.

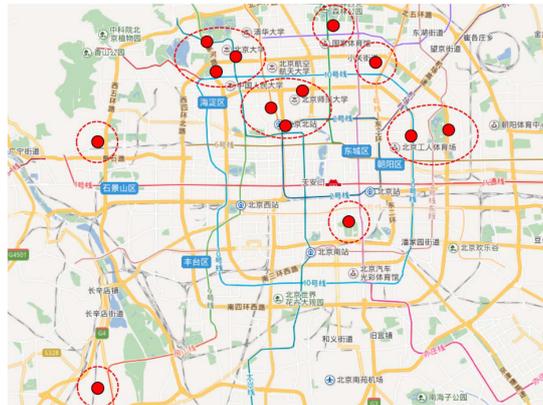
*Discussion* Based on the research questions, we can draw clear conclusions from the findings presented in Table 1. Initially, we compared the performance of different models, including a simple neural network. Although the random forest model outperformed others, its performance was only 3.5 % higher, distinguishing it significantly from models such as KNN and linear regression, which performed 12 % and 13.8 % worse than the random forest, respectively.

The observed performance decrease from 81.5 % to 78.4 % in the FL approach can be attributed to the training process and falls within a reasonable range. Considering the privacy advantages of FL, the 7.2 % accuracy loss compared to the best model is acceptable, particularly when taking into account the reduction in standard deviation from 4.92 to 3.99.

Although this example is hypothetical, it highlights the potential benefits and importance of FL in official statistics. It showcases how FL provides access to crucial data sets for ML while maintaining nearly negligible loss in accuracy compared to a centralized data set.

### 3.2 Fine dust pollution

Reducing air pollution is a significant part of the Sustainable Development Goals (SDGs) established by the United Nations<sup>12</sup>. To measure progress toward achieving SDGs, NGOs and other data producing organizations developed a set of 231 internationally comparable indicators, including *annual mean levels of fine particulate matter (e.g.  $PM_{2.5}$  and  $PM_{10}$ )*. [Hu+18] showed that personalized FL can be used to extract timely high frequent information on air pollution more accurately than models using centralized data.



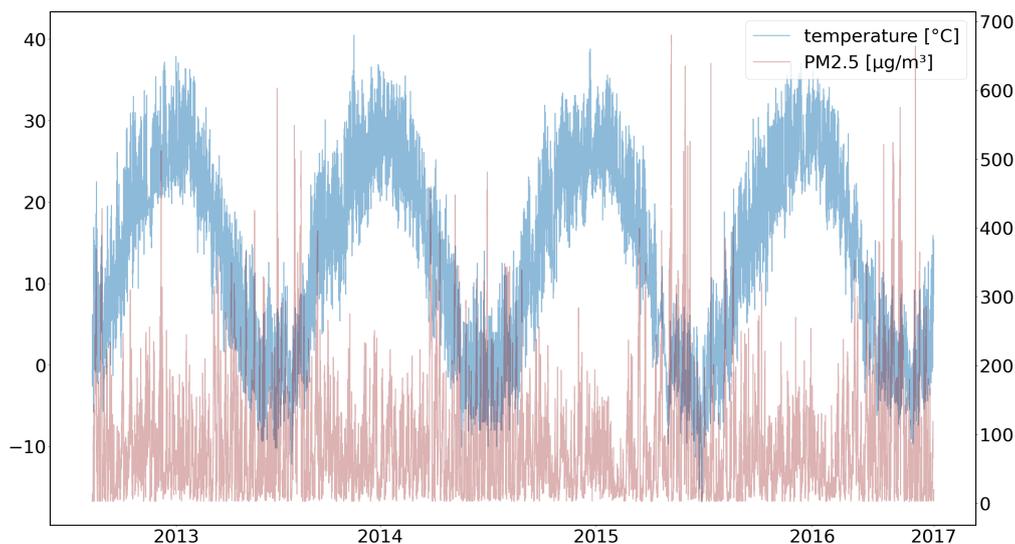
**Fig. 1.** Location of meteorological stations for the fine dust pollution simulation on a map of Beijing, China. 12 of the 13 stations are included in the public data set which we have used for our simulations. The dashed lines mark *regions* of the “Region-Learning” approach in [Hu+18]. Image source: [Hu+18].

In our second use case, we provide a comparison between centralized and FL models (without personalization) and make the developed code and methods accessible. It should be noted that we utilize a slightly different data set and methodology compared to [Hu+18], which we explain at the end of this section. We model a classification task in which the current fine dust pollution is

<sup>12</sup> Air quality and health <https://www.who.int/teams/environment-climate-change-and-health/air-quality-and-health/policy-progress/sustainable-development-goals-air-pollution>, accessed on July 17, 2023

inferred based on meteorological input data. More precisely, 48 consecutive hourly measurements are used to make a prediction for the current  $\text{PM}_{2.5}$  pollution (the total weight of particles smaller than  $2.5\mu\text{m}$  in one  $\text{m}^3$ ). The output of the predictor is one of the three classes *low*, *medium* or *high*. The thresholds for each class are chosen in a way such that the samples of the whole data set are distributed evenly among the three classes.

*Data set* The data set we use is a multi-feature air quality and weather data set [Zha+17] which is publicly available online<sup>13</sup>. It consists of hourly measurements of 12 meteorological stations in Beijing, recorded over a time span of 4 years (2013–2017). Figure 1 depicts the locations of the 12 stations in Beijing. In total, more than 420 000 data records are included in the data set. Although some attributes are missing for some data records, most records have data for all the 17 attributes. An example plot for the two attributes  $\text{PM}_{2.5}$  and temperature is shown in Figure 2.



**Fig. 2.** Example plot for the data of one meteorological station and the two features  $\text{PM}_{2.5}$  and temperature. The four-year time span is clearly visible by the temperature wave, due to hot summers and cold winters.

*Data preprocessing* To complete the missing data records, we use linear interpolation. We apply one-hot encoding to the wind direction attribute. All other features are scaled to a range from 0 to 1. For the attributes  $\text{PM}_{10}$ ,  $\text{SO}_2$ ,  $\text{NO}_2$ ,  $\text{CO}$  and  $\text{O}_3$ , we observe a high correlation with the target attribute and thus exclude them from training. 80% of the data are used as training data, the rest is used as test data.

*Setup* As in the first use case, we implement a centralized learning benchmark and compare it with a FL approach. We model one FL client per meteorological station and split the data accordingly, while the benchmark model is trained with data from all 12 stations. In both settings, we use neural networks with LSTM (long-short term memory) layers and apply 5-fold cross validation. The architecture of the neural networks is similar across both settings and has been manually tuned to reach a good performance: The input layer is followed by a 10-neuron LSTM layer, a dropout layer with a dropout rate of 25%, a 5-neuron LSTM layer, another dropout layer with a dropout rate of 35% and a 3-neuron dense layer for the classification output. For the same reasons as in the first use case, we use the Adam optimizer and apply a learning rate of 0.05 on the server and 0.005 on the client. The client learning rate is decreased every 64 epochs by a factor of 10 to facilitate fine-tuning in later stages of the training. The total training budget we have allocated is 10 epochs for centralized learning and 200 epochs for FL (with a single round of local training per epoch).

<sup>13</sup> Beijing multi-site air-quality data set <https://www.kaggle.com/datasets/sid321axn/beijing-multi-site-airquality-data-set>, accessed on July 17, 2023

*Results* A summary of our results for the fine dust pollution use case is provided in [Table 2](#). Depicted are the means of our 5-fold cross validation experiments.

The centralized learning benchmark reaches a mean classification accuracy of 72.4%, with similarly high numbers for precision and recall (72.8%, respectively 72.3%). In comparison, the FL classifier reaches a performance of both an accuracy and a recall of 68.0% and a precision of 67.9%. The relative standard deviation is higher in the FL scenario for all three metrics, reaching from +2.67 percentage points (accuracy) to +2.9 percentage points (both precision and recall).

An exemplary confusion matrix for one of the five resulting models of the centralized learning is depicted in [Table 3](#). Most misclassifications are made for the *medium* class. The same could be observed for the other models (both in centralized and federated learning).

Model	Accuracy ( $\pm$ std)	Precision ( $\pm$ std)	Recall ( $\pm$ std)	Rel. loss (%)
neural network	72.4% (4.92)	72.8% (8.66)	72.3% (8.10)	0.0
neural network (fed.)	68.0% (7.59)	67.9% (10.05)	68.0% (9.59)	5.9-6.7

**Table 2.** Performance in the fine dust pollution simulation. The span of the relative loss refers to all three metrics.

true class / predicted class	low	medium	high
low	114 450	23 712	5769
medium	23 635	80 718	33 754
high	1944	24 629	111 581

**Table 3.** Exemplary confusion matrix for one of the five models in the cross-validation training of the centralized model for the fine dust pollution use case.

*Discussion* Compared to the first use case, the training database is significantly larger. With 12 clients, there are also four times as many participants in the FL scenario as in the first use case. Still, the performance decrease is small, with an accuracy of 68.0% (FL) compared to 72.4% in the centralized training scenario.

Apart from preprocessing the data set, another time-consuming part of the engineering was tuning the hyperparameters of the FL training. Tools for automatic FL hyperparameter optimization were out of scope for this work, thus it was necessary to manually trigger different trial runs with varying hyperparameters.

*Comparison with literature* The authors of [[Hu+18](#)] compare the results of their personalized FL strategy “Region-Learning” to a centralized learning baseline and standard FL. Although according to the authors, their personalized FL approach outperforms the other two approaches (averaged over the regions by 5 percentage points compared to standard FL), we want to stress that Region-Learning has another goal than standard FL – namely multiple specialized models, and not one global model as in standard FL and most use cases for official statistics (also see [subsection 2.1](#)).

Furthermore, Hu et al. have not provided sufficient information to retrace their experiments. Especially the number of classes for PM<sub>2.5</sub> classification and information on the features used for training the classifiers are missing, so that their results are hard to compare to ours. For example, setting the number of classes to 2 and using all features of the data set (including the other pollution attributes PM<sub>10</sub>, SO<sub>2</sub> etc.) would significantly ease the estimation task. Also, we have no information on whether cross validation was applied in the work of Hu et al. Two more hints in the paper [[Hu+18](#)] suggest that they have used a slightly different data set than we have: The data set they describe includes “more than 100 000” data records from 13 meteorological stations in Beijing, while our data set contains more than 420 000 records from 12 stations.

One consistency across both their work and ours is the accuracy drop from centralized learning to FL, with 4 percentage points in [Hu+18] and 4.4 percentage points in our work.

### 3.3 Mobile radio (LTE)

Mobile Network Operator (MNO) data is a valuable source for obtaining high-frequency and spacial insights in various fields, including population structure, mobility and the socio-economic impact of policy interventions. However, a lack of legal frameworks permitting access to data of all providers, as seen in cases like Germany, constrain the quality of analysis [SBH22]. Accessing only data of selected providers introduces biases, making FL an attractive solution to enhance the representativeness by enabling the aggregation of insights from multiple major MNOs.

Thus, our third use case is based on private MNO data owned by the company umlaut SE<sup>14</sup>. Different from the first two use cases, we had no direct access to the data, just as the aggregation party in realistic FL settings. While this allows for practical insights, it also comes with constricted resources in the private sector. Hence, the focus of this use case is more on practical engineering issues of FL and less on optimal results.

The data set contains mobile communication network coverage data, including latency and speed tests, each linked to the mobile LTE devices of individual users and a specific timestamp. The data records are also associated with GPS coordinates, such that a daily “radius of action” can be computed for each user. This radius describes how far a user has moved from their home base within one day. The user home bases have also been computed on the available data – a home base is defined as the place where most data records have been recorded. The ML task we model in this use case is to estimate the daily radius of action for a user, given different LTE metrics of one particular day (see below).

*Data set* The whole data set originally contains 286 329 137 data records. The following features of the data set have been aggregated for each day and user: *radius of action* in meters, *share of data records with Wi-Fi connection* and the variance and mean values for each of the following LTE metrics: *RSRQ*, *RSRP*, *RSSNR* and *RSSI*. The date has been encoded into three numeric features (*calendar week*, *day of the week* and *month*) and the boolean feature *weekend*.

*Data preprocessing* We set a specific time frame of six months and a geofence around the German state of North Rhine-Westphalia. All other records are excluded – leaving 2 718 416 records in the data set. Additionally, we apply a filtering strategy to clean our data: each user in the database needs to have data for at least 20 different days (within the time span of six months) and 10 records on each of these days. Otherwise, all records of this user are discarded. After the second filtering step, there are 1 508 102 data records in the data set. We scale each feature to a range from 0 to 1 and then use for training, validating and testing our models.

60% of the data are used as training data, 20% are used as validation data and the remaining 20% as test data. For FL, we have divided the data set according to the mobile network operators (MNOs) of the users. Since more than 99.6% of the data records are associated with three major providers, the other 0.4% of the data records (belonging to 29 other MNOs) are eliminated from the data set.

*Setup* We use two centralized learning benchmarks: a random forest regressor and a neural network, which have both been subject to a hyperparameter search prior to their training. The network architecture for both the centralized benchmark neural network and the FL training process is the same: The first layer consists of 28 dense-neurons and the second layer consists of 14 dense-neurons, which lead to the single-neuron output layer. All dense layers except for the output layer use the ReLU activation function. For FL, we use the SGD optimizer with a server learning rate of 3.0, a client learning rate of 0.8 and a batch size of 2.

*Results* The benchmarks of the centralized learning regressors are  $R^2$  values of 0.158 (random forest), 0.13 (neural network) and 0.13 (linear regression). For the neural network trained in the FL scenario, we achieve a slightly lower  $R^2$  value of 0.114 (see Table 4).

<sup>14</sup> umlaut website <https://www.umlaut.com/>, accessed on July 17, 2023

Model	$R^2$	Rel. loss (%)
neural network	0.130	17.7
neural network (federated)	0.114	27.8
random forest	0.158	0.0

**Table 4.** Performance in the mobile radio simulation.

*Discussion* The reasons behind the weak performance of the benchmark models ( $R^2$  of 0.158 and 0.13) are not clear. The hyperparameters might not be optimal, since we were not able to spend many resources on hyperparameter tuning due to time constraints of the data owner. Another reason might be that the modeled task (estimating the radius of action based on LTE connection data) is inherently hard to learn. With an  $R^2$  of 0.114, we were able to reproduce this performance in the FL setting.

Since the private data set in this use case has not left company premises, there are important lessons to be learned from a practical perspective:

1. Even if the data set is not directly available during the model engineering process, it is crucial to get basic insights on the features and statistical properties before starting the training. Essential decisions, such as the type of model to be trained, can be made based on this.
2. A thorough hyperparameter optimization is needed to obtain useful results. It might take a lot of time and computational resources to find hyperparameters which are suited for the task.
3. Technical difficulties while creating the necessary APIs and setting up the chosen ML framework at the FL clients can slow down the process even more. Without access to the database, it might be hard to reproduce technical errors.

While all points mentioned above were encountered in the third simulation, there was only *one* party who held all data. In real FL scenarios with multiple data holders, the process might get much more complicated.

## 4 Key Observations

Our simulations lead to the following key observations:

Models trained via FL can reach a performance very close to models trained with centralized ML approaches, as we have shown in all three use cases. While the performance gap itself is not surprising (since the FL model has been exposed to the complete data set only indirectly), we want to stress that without FL, many ML scenarios might not be possible due to privacy concerns, trade secrets, or similar reasons. This is especially true for health care data, i.e., the domain of our first simulation.

While the random forest regressor has demonstrated superior performance compared to other centralized learning benchmarks in all three simulations, exploring the potential of tree-based models within a FL context [AI+23; YOW22; LWH20] could be a promising avenue for further investigation. The improved interpretability and explainability over many other models, e.g., neural networks, is another advantage of tree-based models.

On the other hand, random forest regressors are not suitable if tasks get more complicated. Also, their architecture, i.e., many decision trees which may be individually overfitted to parts of the training data, can facilitate the extraction of sensitive information of the training data and thus pose an additional privacy risk.

Choosing the right hyperparameters is crucial for any ML model. Since automatic HPO is still an open problem for FL algorithms, (manually) finding the right settings can be a time-consuming process. Developing a suitable framework for automated HPO for FL would be important future work – although for official statistics, other issues might be more pressing at the moment (see [section 5](#)).

In our third simulation (mobile radio data), we did not have access to the training and test data set, just like in a real-world scenario. This means both HPO and technical debugging needed to be performed remotely, without access to the data. Although this was already challenging, we

believe that in scenarios with multiple data holders and possibly heterogeneous data sets, these tasks will be even harder.

All FL simulations were performed on the machine which also had access to the complete data set. In a real-world application, where each client runs on a distinct machine, other settings and other frameworks might be more practical than TensorFlow Federated.

Last but not least, we want to emphasize that FL, despite its privacy-enhancing character, may still be vulnerable to some ML privacy issues (see [subsection 2.2](#)). Hence, analyzing and communicating these risks is an important step before an application is rolled out in practice.

## 5 Implications for Official Statistics

In this work, we have demonstrated how FL can enable NSOs to address pressing data needs in fields that are relevant to policymakers and society. Official statistics are characterized by high accuracy while underlying strict standards in confidentiality and privacy. Accuracy, explainability, reproducibility, timeliness, and cost-effectiveness are essential quality dimensions for statistical algorithms [Yun+22]. In this setting, our findings indicate that FL bears significant potential to support statistical production and improve data quality.

We have shown that FL can empower NSOs to generate reliable models that accurately capture global relations. In each of our use cases, the FL-generated models exhibited nearly identical predictive performance compared to a model created by combining all available data. Each model architecture that performed well on centralized or local data could be easily adapted to a FL training process with a similar level of predictive performance only using distributed data.

If upcoming applications require to optimize an individual model for each participating party, personalized FL can be used to generate potentially improved models tailored to individual clients. This increases the interest to cooperate for each participating party, as it offers to enhance the analytic potential for each client and the server. However, it is important to note that this customization may come at the cost of global predictive performance.

FL provides the main advantage of not needing to exchange sensitive data (see [subsection 2.2](#)). Additionally, there is no need to store or process the complete data set centralized in the NSOs.

NSOs can be empowered to appraise novel data sources sans the need for new legislation. In cases where legislative changes prove impractical, FL provides a crucial pathway to assess and prepare for regulations' modernization. By showcasing the advantages and implications of accessing new data sources before legal frameworks permit, FL not only significantly accelerates and relieves statistics production but also occasionally enables it.

To ensure successful future implementations of FL in NSOs, it is essential to focus on further advancements. Specifically, improvements in communication frequency are crucial to enable high-speed and efficient exchanges. Our observations indicate that FL generally requires a greater number of epochs (distributed across communication rounds) compared to centralized training to achieve similar performance levels. In our use cases, even with small datasets, we found that at least 50 rounds of communication were necessary. In real-world applications, this would result in high delay and cost. Therefore, the development of infrastructure for seamless sending and receiving ML models is necessary. Addressing this challenge, we discovered that the implementation of adaptive server optimization techniques reduced the training rounds and contributed to training stability. As a result, we recommend the use of adaptive optimizers to help minimize communication costs and enhance the efficiency of FL processes. By incorporating such adaptive optimization methods, NSOs can optimize the performance and effectiveness of FL while reducing the burden of communication overhead.

Additionally, it is crucial to provide partners with the necessary tools to update models effectively. This requires coordination of the server and expertise from all participating parties. In practice, real-world applications of FL often involve the challenge of harmonizing client data without directly accessing it. Achieving an optimized model architecture uniformly across all clients also necessitates the knowledge and collaborative efforts of the clients themselves. Providing comprehensive tools and resources to partners enables them to actively contribute to the model updating process while maintaining data privacy and security.

FL is evolving rapidly and both industry and research will continue to improve the field in the coming years. The performance and efficiency of practical FL frameworks is expected to be further optimized. Similarly, we expect the development of more usable PPML algorithms including the ones based on Secure Multi-Party Computation (SMPC) and Homomorphic Encryption (HE) –

allowing for provably secure collaborative ML. Although such PPML methods have been proposed and frameworks exist, their performance today is often far from acceptable for many practical applications. With more standardization and simpler, respectively more efficient, applications, FL will become even more beneficial to official statistics.

In summary, FL should indeed be recognized as an important technology that can facilitate the modernization of legal frameworks for official statistics. It enables NSOs to safely use publicly relevant information that is not expected to be accessed by future legal frameworks, ultimately enhancing the quality and relevance of official statistics. However, further development is still required to fully realize the potential of FL in this context.

## 6 Conclusion

In scenarios where external partners are unwilling to share individual-level information but still aim to analyze or disseminate global insights in their field of application, FL can help to overcome these issues. We have shown across a range of three simulated use cases that FL can reach a very similar performance to centralized learning algorithms. Hence, our results indicate that if classic (centralized) ML techniques work sufficiently well, FL can possibly produce models with a similar performance.

One of the next steps to transfer FL into the practice of official statistics could be to conduct practical pilot studies. These could further showcase both the applicability and challenges of FL beyond a simulated context. Another focus of future work in this area could be the analysis of privacy risks in FL scenarios of official statistics and potential mitigation strategies. This would be an important stepping stone in ensuring the privacy protection of involved parties, on top of the privacy enhancement by using FL. Just as in countless other domains, we expect FL to become a relevant technology for official statistics in the near future.

## References

- [Al+23] Mohammad Al-Quraan et al. *FedTrees: A Novel Computation-Communication Efficient Federated Learning Framework Investigated in Smart Grids*. 2023. DOI: [10.1016/j.engappai.2023.106654](https://doi.org/10.1016/j.engappai.2023.106654).
- [BDF18] Martin Beck, Florian Dumpert, and Joerg Feuerhake. “Machine learning in official statistics”. In: *arXiv preprint arXiv:1812.10422* (2018). URL: <https://arxiv.org/abs/1812.10422>.
- [Buc23] David Buckley. *15. United Nations Economic Commission for Europe: Trialling approaches to privacy-preserving federated machine learning*. 2023. URL: <https://unstats.un.org/wiki/display/UGTTOPPT/15.+United+Nations+Economic+Commission+for+Europe%3A+Trialling+approaches+to+privacy-preserving+federated+machine+learning>.
- [DB17] Florian Dumpert and Martin Beck. “Einsatz von Machine-Learning-Verfahren in amtlichen Unternehmensstatistiken”. In: *AStA Wirtschafts-und Sozialstatistisches Archiv* 2.11 (2017), pp. 83–106.
- [Eur22] United Nations Economic Commission for Europe. “Machine Learning for Official Statistics”. In: *UNECE Machine Learning Group* (2022). URL: <https://unece.org/statistics/publications/machine-learning-official-statistics>.
- [HAP17] Briland Hitaj, Giuseppe Ateniese, and Fernando Perez-Cruz. “Deep models under the GAN: information leakage from collaborative deep learning”. In: *ACM CCS*. 2017, pp. 603–618.
- [Hu+18] Binxuan Hu et al. “Federated region-learning: An edge computing based framework for urban environment sensing”. In: *IEEE GLOBECOM*. IEEE. 2018, pp. 1–7.
- [KKP20] Viraj Kulkarni, Milind Kulkarni, and Aniruddha Pant. “Survey of personalization techniques for federated learning”. In: *2020 Fourth World Conference on Smart Trends in Systems, Security and Sustainability (WorldS4)*. IEEE. 2020, pp. 794–797.
- [LWH20] Qinbin Li, Zeyi Wen, and Bingsheng He. “Practical Federated Gradient Boosting Decision Trees”. In: *Proceedings of the AAAI Conference on Artificial Intelligence* 34.04 (2020), pp. 4642–4649. URL: <https://doi.org/10.1609/aaai.v34i04.5895>.

- [McM+17] Brendan McMahan et al. “Communication-efficient learning of deep networks from decentralized data”. In: *Artificial intelligence and statistics*. PMLR. 2017, pp. 1273–1282.
- [Mel+19] Luca Melis et al. “Exploiting unintended feature leakage in collaborative learning”. In: *2019 IEEE symposium on security and privacy (SP)*. IEEE. 2019, pp. 691–706.
- [NSH19] Milad Nasr, Reza Shokri, and Amir Houmansadr. “Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning”. In: *IEEE SP*. IEEE. 2019, pp. 739–753.
- [Red+20] Sashank Reddi et al. “Adaptive federated optimization”. In: (2020). DOI: [ARXIV. 2003.00295](https://arxiv.org/abs/2003.00295). URL: <https://arxiv.org/abs/2003.00295>.
- [SBH22] Younes Saidani, Sarah Bohnensteffen, and Sandra Hadam. “Qualität von Mobilfunkdaten – Projekterfahrungen und Anwendungsfälle aus der amtlichen Statistik”. In: *WISTA - Wirtschaft und Statistik* 5 (2022), pp. 55–67. URL: <https://www.destatis.de/DE/Methoden/WISTA-Wirtschaft-und-Statistik/2022/05/qualitaet-mobilfunkdaten-052022.html>.
- [Sho+17] Reza Shokri et al. “Membership inference attacks against machine learning models”. In: *2017 IEEE symposium on security and privacy (SP)*. IEEE. 2017, pp. 3–18.
- [SRS17] Congzheng Song, Thomas Ristenpart, and Vitaly Shmatikov. “Machine learning models that remember too much”. In: *Proceedings of the 2017 ACM SIGSAC Conference on computer and communications security*. 2017, pp. 587–601.
- [Tem22] Julian Templeton. *Privacy enhancing technologies: An overview of federated learning*. 2022. URL: <https://www.statcan.gc.ca/en/data-science/network/privacy-enhancing-techniques> (visited on 09/29/2023).
- [YOW22] Fuki Yamamoto, Seiichi Ozawa, and Lihua Wang. “eFL-Boost: Efficient Federated Learning for Gradient Boosting Decision Trees”. In: *IEEE Access* 10 (2022), pp. 43954–43963. DOI: [10.1109/ACCESS.2022.3169502](https://doi.org/10.1109/ACCESS.2022.3169502).
- [Yun+22] Wesley Yung et al. “A quality framework for statistical algorithms”. In: *Statistical Journal of the IAOS* 38.1 (2022), pp. 291–308.
- [YZH21] Xuefei Yin, Yanming Zhu, and Jiankun Hu. “A comprehensive survey of privacy-preserving federated learning: A taxonomy, review, and future directions”. In: *ACM Computing Surveys (CSUR)* 54.6 (2021), pp. 1–36.
- [Zha+17] Shuyi Zhang et al. “Cautionary tales on air-quality improvement in Beijing”. In: *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* 473.2205 (2017), p. 20170457.
- [ZLH19] Ligeng Zhu, Zhijian Liu, and Song Han. “Deep leakage from gradients”. In: *Advances in neural information processing systems* 32 (2019).