

Empirical Analysis of Bias in Federated Learning

Nawel Benarba, Sara Bouchenak

INSA Lyon – LIRIS

Lyon, France

nawel.benarba@insa-lyon.fr, sara.bouchenak@insa-lyon.fr

Résumé

Federated Learning (FL) is a machine learning paradigm that allows distributed clients to collaboratively train a global model without having to share their local data, preserving data privacy. However, it presents new challenges, including the potential for models to exhibit bias towards specific demographic groups. Motivated by this inherent issue, we conduct an extensive empirical analysis where we measure FL bias through the disparity in terms of model quality and demographic parity. First, we conduct an empirical evaluation on four widely used datasets, to evaluate the impact of data size and heterogeneity on FL model bias. Then, we analyze the actual effectiveness of the state-of-the-art bias mitigation methods on different datasets. Our findings reveal interesting observations indicating that an increase on data size or heterogeneity level comes with an increase of FL bias. It also shows that the bias mitigation mechanisms are more effective for datasets with less FL bias.

Mots-clés : Federated Learning, Bias, Empirical Analysis

1. Introduction

Computer technology growth has made huge data accessible. However, centralizing this data comes with high costs, lengthy processes and privacy concerns. Traditional centralized machine learning (ML) cannot handle the increasing complexity and diversity of data, which makes growing interests to federated learning. A distributed ML paradigm that enables different clients to collaboratively train a global model without sharing their local data. It has received significant attention for its ability to train large-scale models in a decentralized manner while providing user data privacy. It has been applied in crucial areas, such as healthcare [2], loan granting [19], IoT [18], etc. FL applications are increasingly facing bias, resulting in unfair decisions that exhibit discriminatory behavior towards certain demographic groups based on their sensitive attributes, including gender, race, and age [3]. One notorious instance is the COMPAS algorithm [5], which predicts the likelihood of a criminal reoffending. The algorithm demonstrate significant racial bias due to the use of historical biased training records. In our analysis we assess FL bias through a series of experiments, we outline our primary contributions below. We consider two FL bias aspect measuring model quality and demographic parity disparities on four datasets. We show how does data heterogeneity impacts bias. We investigate how the amount of training data affects bias. We evaluate state-of-the-art bias mitigation. Finally, we discuss and identify novel research directions for FL bias mitigation.

Our results show that larger or more heterogeneous datasets increase the potential for FL models to exhibit bias. It also indicates that bias mitigation mechanisms are more effective with datasets that initially exhibit low FL bias using demographic parity metric. Based on our analysis, we identify several research areas. For instance, exploring the compatibility of bias mitigation methods with different metrics could lead to more effective solutions. Additionally, reducing the size of training data could be a viable approach to mitigate bias.

Paper Roadmap. The rest of the paper is structured as follows. First, we provide the necessary background on FL and bias in Section 2. Our methodology is presented in Section 3. Our extensive experimental results are presented in Section 4. In Section 5, we discuss our findings and identify open research directions. We survey the related work in Section 6, and conclude in Section 7.

2. Background

2.1. Federated Learning

FL is a learning paradigm in which a server orchestrates different distributed clients (c_1, c_2, \dots, c_N) and aggregates their local models θ_i to obtain a single global model θ using massive amounts of data kept private on each device. The key idea is to minimize the expected value of a loss function using a dataset $d = (d_1, d_2, \dots, d_n) \in D$ of n samples distributed across k involved clients. To do so, the federated averaging algorithm, FedAvg [16], samples a subset of K clients per round to perform local training of the global model on their local datasets. The model updates are then averaged at the server, being weighted based on the size of their respective datasets according to Equation 1.

$$\theta = \sum_{i=1}^{i=k} \frac{n_i}{|D|} \cdot \theta_i \quad (1)$$

2.2. Bias in Federated Learning

Despite the potential benefits of FL in training machine learning models on decentralized data, it is not immune to bias which appears due to the use of incomplete, faulty or prejudicial datasets and models, producing unfair models towards certain groups based on their sensitive attributes. The disparity between the predictions on the privileged group and those on the unprivileged group is used to quantify bias. A variety of metrics have been proposed to measure it [17]. In our empirical analysis we consider both of Statistical Parity Difference and Equal Opportunity Difference.

Statistical Parity Difference (SPD). Also known as demographic parity, reflects the property that the demographics of those receiving positive classifications are identical to the demographics of the population as a whole [9], a binary predictor \hat{Y} is fair from the statistical parity perspective with regard to the set of the considered sensitive attributes $S_A = \{S_1, S_2, \dots, S_{s_a}\}$, if for each sensitive attribute $S_j \in S_A$ we have $\Pr(\hat{Y} = 1|S_j = 1) = \Pr(\hat{Y} = 1|S_j = 0)$. The Statistical Parity Difference (SPD) metric is defined as :

$$SPD_{S_j} = \Pr(\hat{Y} = 1|S_j = 0) - \Pr(\hat{Y} = 1|S_j = 1) \quad (2)$$

Equal Opportunity Difference (EOD). Unlike the statistical parity difference, equal opportunity difference allows the binary predictor \hat{Y} to depend on S_j but only through the target variable Y , imposing a constraint on bias as well as on model quality. It requires non-discrimination within the advantage outcome group [12]. The predictor is considered fair if the true positive rate is independent of the sensitive attributes. To measure this, EOD metric is calculated as :

$$EOD_{S_j} = \Pr(\hat{Y} = 1|S_j = 0, Y = 1) - \Pr(\hat{Y} = 1|S_j = 1, Y = 1) \quad (3)$$

3. Methodology

To investigate the presence of bias in FL and assess how it can be impacted, we quantify it through two metrics, SPD and EOD, with four widely used datasets. Our analysis follows an experimental design approach [11], guided by a set of research questions :

- What is the presence of FL bias in different datasets, and with regard to different sensitive attributes ?
- How does data heterogeneity impact bias ?
- To what extent do data scale influence bias ?
- What is the impact of state-of-the-art FL bias mitigation mechanisms on different datasets ?

Different Bias Metrics. To accurately and comprehensively measure bias in FL, we select the commonly used SPD and EOD metrics to cover two different bias aspect. SPD measures the disparity in outcomes among different demographic groups, while EOD focuses on the precision of the model's predictions. These metrics are intended for binary sensitive attributes. For the non-binary ones, we convert them into binary form by applying a hierarchical clustering algorithm.

Different Datasets. To carry out our investigation and identify discernible patterns or trends that emerge, we opt to use four real-world datasets, that are recognized for their inherent bias issues in the context of FL applications.

Bias with Regard to Different Sensitive Attributes. To conduct a thorough evaluation of FL bias, we consider the different sensitive attributes present within each dataset. This approach allows us to evaluate potential sources of bias across a range of datasets with regard to the differences in sensitive attributes that may exist between them.

FL Data Heterogeneity. FL data heterogeneity due to users demographics refers to the variation in demographic characteristics of datasets collected from different nodes in a FL system. Since our objective is to analyze the influence of data heterogeneity on FL bias, we explore various scenarios by varying Dirichlet parameter [13] to manipulate data heterogeneity level. Using the resulting data distributions, we assess bias at the model level.

FL Data Size. With the aim of evaluating the impact of data size on FL bias, we investigate different scenarios by adjusting the amount of data used by FL clients. We gradually increase the training distributions size in each scenario and measure bias using SPD and EOD metrics.

FL Bias Mitigation. To assess the effectiveness of state-of-the-art bias mitigation methods, we use FCFL as it is compatible with the metrics used in our analysis [22]. We compare it with the conventional FL aggregation method to determine the reduction of bias.

4. Empirical Analysis of Bias in Federated Learning

4.1. Experimental Setup

In this section, we perform a comprehensive analysis of bias in FL across various system setups. Our objective is to understand how the magnitude of FL bias evolves with respect to the scale of data, and the levels of heterogeneity in data distributions among clients. We also examine the performance of reducing bias by contrasting FCFL with the conventional FL approach. We partition the datasets into a training set (80%), a validation set (10%), and a test set (10%). We note that the default FL aggregation algorithm we use is FedAvg with non-IID data generated by Dirichlet function with an alpha parameter of 0.02. We consider different sensitive attributes and FL clients number depending on the dataset as follows : for KDD and Adult, we consider age as sensitive attribute with respectively 5 and 10 clients ; for DC and MEPS, sensitive attributes are respectively gender with 4 clients for the former, and race with 10 clients for the latter.

We conduct our experiments on a testbed platform that consists of two Intel Xeon E5-2680 v4 processors with 14 cores each, 128 GB of RAM, and two NVIDIA GeForce GTX 1080 Ti GPUs. For the implementation, we use Python 3.9 and PyTorch 1.10 in all of our experiments, with the exception of FCFL, where we rely on the official repository's implementation and environment [6], which employed Pytorch 1.6 and Python 3.7.

4.2. Datasets and Models

To evaluate bias in FL and the factors that impact it, we use four real-world datasets, namely Adult, KDD, DC, and MEPS, which are renowned for their inherent bias issues.

KDD. KDD dataset contains 299,285 records derived from weighted census data obtained from the 1994 and 1995 population surveys in the United States [8]. The dataset's objective is to predict whether an individual earns over 50k \$ per year, the sensitive attributes are age, gender, and race.

DC. The dataset consists of 60,420 records from census information in the Netherlands in 2001 [20]. Its objective is to classify an individual's occupation as having either high or low prestige, with sensitive attributes being age and gender.

MEPS. A collection of healthcare data and expenditure information from a nationally representative sample of US households [4]. It contains 35,428 records, including gender and race as sensitive attributes. It is used to predict whether an individual has used medical facilities more than 10 times.

Adult. A widely used dataset for evaluating bias in ML models [14]. It comprises 48,842 records, its prediction task and sensitive attributes are the same as KDD.

To obtain valuable insights and accurate predictions from the four datasets, we use logistic regression. Each FL system is initialized with the same initial model through transfer learning to ensure an impartial evaluation of different FL systems. This initial model provides a common starting point for setting a performance baseline, which can be used to assess the effectiveness of subsequent models.

4.3. Evaluation of Bias in Federated Learning

In order to assess bias in federated learning, we quantify both of SPD and EOD metrics for each dataset under the default configuration. The selection of the sensitive attributes varies with age, gender, and race being considered for KDD and Adult, gender and age for DC, and gender and race for MEPS.

Figure 1 presents SPD and EOD values (y axis) for each dataset (x axis) with regard to its different sensitive attributes. The results reveal that both SPD and EOD metrics effectively capture bias and exhibit similar trends across various sensitive attributes. However, the EOD metric consistently yields lower values than SPD according to the model quality, we refer this to the EOD's additional constraint. It measures the difference between the correct positive predictions, which contributes to its lower values compared to SPD that measures the difference between all the positive predictions.

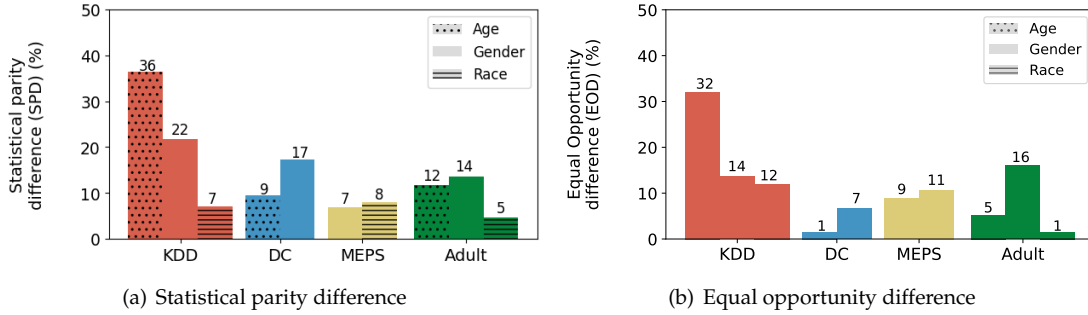


FIGURE 1 – Bias in FL on different datasets, with regard to different sensitive attributes

4.4. Heterogeneity and Bias in Federated Learning

In this section, we aim to investigate how data heterogeneity affects bias in the context of FL. To achieve this, we generate non-identical clients using Dirichlet process, varying α parameter that controls how identical the clients are. Smaller α values result in a more heterogeneous distribution of data across clients, while larger values lead to more similar and uniform data. We evaluate three different levels of heterogeneity : $\alpha = 0.1$, $\alpha = 0.5$, and $\alpha = 1$, using KDD dataset with age as sensitive attribute, since it exhibits the highest level of bias for the selected dataset (see Figure 1).

Figure 2 shows the SPD and EOD values for the considered Dirichlet parameters, each row corresponds to bias with regard to a particular sensitive attribute. We observe that increasing the degree of heterogeneity with regard to age increases bias with regard to this sensitive attribute while it remains relatively stable with regard to the other attributes. Our findings indicate that an increase in heterogeneity level with regard to a particular sensitive attribute introduces more bias towards that attribute and results in a corresponding rise in SPD and EOD metrics, indicating more bias.

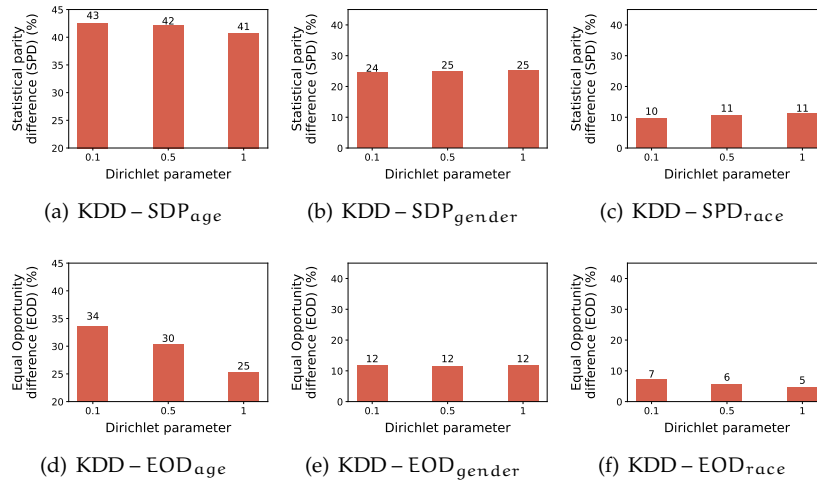


FIGURE 2 – Bias with different levels of FL data heterogeneity with regard to age sensitive attribute

4.5. Federated Learning Bias at Scale

To understand how bias is affected by the amount of data used in the learning process, we perform a series of experiments on KDD dataset split in several partitions, each containing a different percentage of data records of the whole dataset. We progressively increase the amount of used data. Starting from an initial training distribution of 10% to 100%, with a 10% increment. We include the data from each partition in the ones with a larger number of records. Namely, the 10% partition is included in the 20% partition, and so on, until reaching the 100% partition, which contains the whole KDD train set. This gradual expansion allows us to expose the model's learning process to progressively larger and more diverse data distributions. We analyze the corresponding results of SPD and EOD, presented throughout the amount of data in Figure 3.

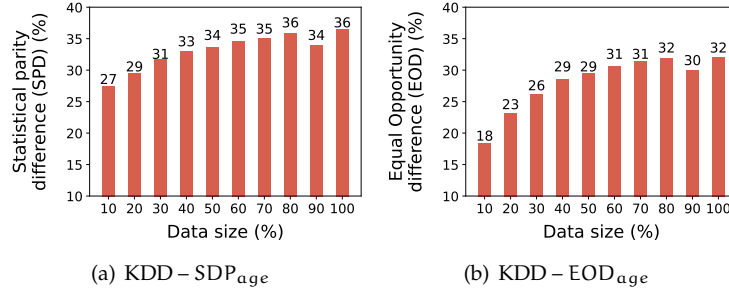


FIGURE 3 – Bias with different data sizes

The results show that bias metrics tend to increase with the expansion of the training data size. This trend can be explained by the fact that larger data size can amplify the potential for the presence of bias in the learning process. As the training data size increases, the model may be more prone to picking up spurious correlations or patterns that reflect bias. Furthermore, the model accuracy remains consistent across the different data distributions as shown in Figure 4, which depicts the model accuracy for the different data sizes. This indicates that the increase of bias is not due to a decline in model performance but rather to the emergence of bias in the learning process.

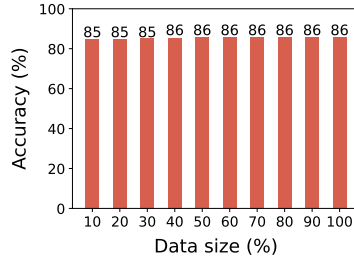


FIGURE 4 – Model accuracy at scale – KDD dataset

4.6. Evaluation of Federated Learning Bias Mitigation

To evaluate the effectiveness of bias mitigation in FL, we compare FCFL approach to the conventional FL aggregation method, Vanilla, which has no bias mitigation policy. We conduct our experiments under the default configuration and a bias constraint of 3% for both SPD and EOD metrics.

Figure 5 shows the evolution of SPD and EOD metrics assessing bias of the aggregated global model throughout FL rounds. We notice that the use of SPD as a bias metric allows the approach to successfully mitigate bias for the datasets that present less bias, namely Adult and MEPS. However, for the datasets with more bias, KDD and DC, FCFL fails to meet the bias constraint. This can be attributed to the fact that the pareto solution to which the method converges does not result in global bias mitigation. Furthermore, the use of EOD shows that FCFL reduces disparities without meeting the defined threshold, indicating that it is still effective in mitigating bias, but not to the required extent. Overall, it is important to consider the type and amount of bias present in the dataset when selecting a bias metric and evaluating the effectiveness of a bias mitigation approach.

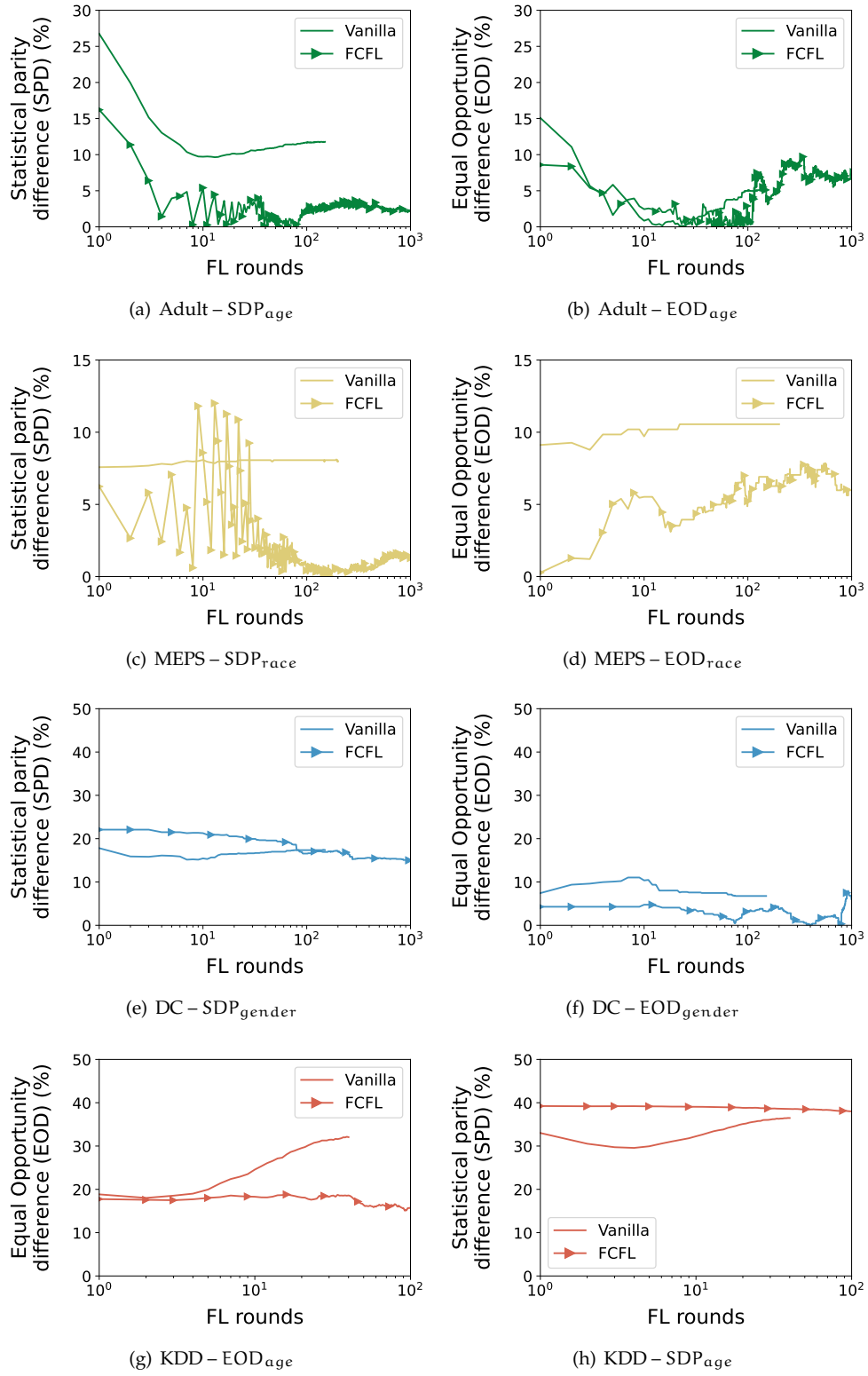


FIGURE 5 – Impact of state-of-the-art FL bias mitigation

5. Discussion and Research Directions

Incorporating SPD and EOD metrics into our analysis provides a deeper understanding of the behavior of bias in FL. Both of the used metrics have the same variation trend. However, EOD captures less bias than SPD due to its additional constraint. The results of our empirical analysis show also that increasing the level of data heterogeneity with regard to a particular sensitive attribute introduces more bias towards that attribute. Similarly, increasing the amount of data in the learning process raises the potential for bias. To address this issue, we found that FCFL is an effective bias mitigation mechanism, particularly with datasets that initially exhibit low levels of bias using SPD metric. Nevertheless, achieving the bias constraint in general remains challenging. The analysis allows us to identify several areas that require further investigation and exploration including :

Leveraging Statistical Data Distributions for Bias Reduction. One promising direction for future research is to build novel FL protocols that gather and are aware of statistical information about FL clients' data distributions. This information can be used to adjust the model update process. By doing so, it may be possible to mitigate the FL bias in the resulting models.

Reducing the Impact of Data Size on Bias. Another research direction is to determine the required amount of data needed in the learning process to ensure an accurate model without increasing bias. This would involve determining the optimal proportion of data each client should use based on its bias.

Achieving More Efficient and Fair Solutions by Mitigating Several Bias Metrics Simultaneously. It is important to assess the compatibility of bias mitigation mechanisms with various bias metrics and develop a bias mitigation framework that can tackle multiple non-controversial bias metrics simultaneously. Such mechanism could provide more efficient and effective solutions.

6. Related work

Bias is a significant concern in machine learning that occurs due to biased or unrepresentative training data, skewed sampling, and user-based biases. Similarly, federated learning is vulnerable to various biases that have the potential to negatively impact model performance and result in unfair outcomes. In our analysis, we evaluate the presence of bias in FL and explore how it can be impacted. Several methods were proposed to detect and assess the potential biases in FL. These methods can be grouped into two main categories : data-level evaluation and model-level evaluation. Data-level evaluation focuses on evaluating the distribution of the training data. For instance, Wasserstein distance [21] is used to measure the degree of imbalance and inequality in data distributions across FL clients. On the other hand, model-level evaluation methods have been less explored compared to the data-level ones. These methods focus on evaluating the performance of the model trained through FL. Several metrics have been proposed to measure the bias present in model predictions such as group fairness metrics [10] and the federated version of Gini coefficient [15]. Our experimental analysis falls under the latter category, assessing bias of the model across different demographic groups. Recent research in FL bias evaluation has also focused on developing novel techniques that seeks to ensure that FL models are fair. Pre-processing techniques such as local data reweighting [1] offer local bias mitigation, clients data samples are reweighted according to their sensitive groups. However, due to the data heterogeneity between FL clients it may not provide global bias mitigation. In-processing techniques are currently the most explored techniques [1, 7, 9, 22]. Among these techniques, FCFL [22] seeks to solve a multi-objective constrained optimization problem, where the bias constraint is enforced for each client, while maintaining consistent model quality.

7. Conclusion

Bias is a significant concern in federated learning, as it can lead to models that are less effective for certain demographic groups, perpetuating inequalities and potentially having serious consequences. To address this issue, it is important to evaluate and monitor bias. In our empirical analysis, we measure bias with two distinct metrics, while exploring how it is affected by data size and heterogeneity levels. Furthermore, we examine how it can be reduced by existing bias mitigation mechanisms. We identify several research avenues that could be pursued, addressing data size and heterogeneity impact to mitigate bias, exploring the compatibility of bias mitigation mechanisms with various bias metrics to develop a framework that mitigates multiple metrics simultaneously for more efficient solutions.

Bibliographie

1. Abay (A.), Zhou (Y.), Baracaldo (N.), Rajamoni (S.), Chuba (E.) et Ludwig (H.). – Mitigating Bias in Federated Learning, 2020.
2. Antunes (R. S.), da Costa (C. A.), Küderle (A.), Yari (I. A.) et Eskofier (B. M.). – Federated Learning for Healthcare : Systematic Review and Architecture Proposal. *ACM Transactions on Intelligent Systems and Technology*, vol. 13, n4, 2022, pp. 54 :1–54 :23.
3. Chen (J.), Kallus (N.), Mao (X.), Svacha (G.) et Udell (M.). – Fairness Under Unawareness : Assessing Disparity when Protected Class is Unobserved. – In danah boyd et Morgenstern (J. H.) (édité par), *Proceedings of the Conference on Fairness, Accountability, and Transparency, FAT* 2019, Atlanta, GA, USA, January 29-31, 2019*, pp. 339–348, Atlanta, GA, USA, 2019.
4. Cohen (S.). – Design Strategies and Innovations in the Medical Expenditure Panel Survey. *Medical care*, vol. 41, 08 2003, pp. III5–III12.
5. Corbett-Davies (S.), Pierson (E.), Feller (A.) et Goel (S.). – A Computer Program Used for Bail and Sentencing Decisions was Labeled Biased Against Blacks. It's Actually not that Clear. *Washington Post*, vol. 17, 2016.
6. Cui (S.), Pan (W.), Liang (J.), Zhang (C.) et Wang (F.). – Repository of FCFL Software for Bias Mitigation in Federated Learning, 2021.
7. Du (W.), Xu (D.), Wu (X.) et Tong (H.). – Fairness-aware Agnostic Federated Learning. – In Demiconi (C.) et Davidson (I.) (édité par), *Proceedings of the 2021 SIAM International Conference on Data Mining, SDM 2021, Virtual Event, April 29 - May 1, 2021*, pp. 181–189, , 2021.
8. Dua (D.) et Graff (C.). – UCI Machine Learning Repository, 2017.
9. Ezzeldin (Y. H.), Yan (S.), He (C.), Ferrara (E.) et Avestimehr (S.). – FairFed : Enabling Group Fairness in Federated Learning, 2021.
10. Garg (P.), Villasenor (J. D.) et Foggo (V.). – Fairness Metrics : A Comparative Analysis. – In Wu (X.), Jermaine (C.), Xiong (L.), Hu (X.), Kotevska (O.), Lu (S.), Xu (W.), Aluru (S.), Zhai (C.), Al-Masri (E.), Chen (Z.) et Saltz (J.) (édité par), *2020 IEEE International Conference on Big Data (IEEE BigData 2020), Atlanta, GA, USA, December 10-13, 2020*, pp. 3662–3666, Atlanta, GA, USA, 2020.
11. Gregar (J.). – Research Design (Qualitative, Quantitative and Mixed Methods Approaches). *Book published by SAGE Publications*, vol. 228, 1994, p. .
12. Hardt (M.), Price (E.) et Srebro (N.). – Equality of Opportunity in Supervised Learning. – In Lee (D. D.), Sugiyama (M.), von Luxburg (U.), Guyon (I.) et Garnett (R.) (édité par), *Advances in Neural Information Processing Systems 29 : Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pp. 3315–3323, Barcelona, Spain, 2016.
13. Hsu (T. H.), Qi (H.) et Brown (M.). – Measuring the Effects of Non-Identical Data Distribution for Federated Visual Classification, 2019.
14. Kohavi (R.). – Scaling Up the Accuracy of Naive-Bayes Classifiers : A Decision-Tree Hybrid. – In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining, KDD'96, KDD'96*, p. 202–207, Portland, Oregon, USA, 1996.
15. Li (X.), Zhao (S.), Chen (C.) et Zheng (Z.). – Heterogeneity-aware Fair Federated Learning. *Information Sciences*, vol. 619, 2023, pp. 968–986.
16. McMahan (B.), Moore (E.), Ramage (D.), Hampson (S.) et y Arcas (B. A.). – Communication-Efficient Learning of Deep Networks from Decentralized Data. – In Singh (A.) et Zhu (X. J.) (édité par), *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, AISTATS 2017, 20-22 April 2017, Fort Lauderdale, FL, USA, Proceedings of Machine Learning Research*, volume 54, pp. 1273–1282, Fort Lauderdale, FL, USA, 2017.
17. Mehrabi (N.), Morstatter (F.), Saxena (N.), Lerman (K.) et Galstyan (A.). – A survey on bias and fairness in machine learning. *ACM Computing Surveys*, vol. 54, n6, 2022, pp. 115 :1–115 :35.
18. Nguyen (D. C.), Ding (M.), Pathirana (P. N.), Seneviratne (A.), Li (J.) et Poor (H. V.). – Federated Learning for Internet of Things : A Comprehensive Survey. *IEEE Communications Surveys and Tutorials*, vol. 23, n3, 2021, pp. 1622–1658.
19. Shingi (G.). – A Federated Learning Based Approach for Loan Defaults Prediction. – In Fatta (G. D.), Sheng (V. S.), Cuzzocrea (A.), Zaniolo (C.) et Wu (X.) (édité par), *20th International Conference on Data Mining Workshops, ICDM Workshops*, pp. 362–368, Sorrento, Italy, 2020.
20. Van der Laan (P.). – The 2001 Census in The Netherlands. – In *Conference the Census of Population*,

Washington, USA, 2000.

21. Zhang (F.), Kuang (K.), Liu (Y.), Wu (C.), Wu (F.), Lu (J.), Shao (Y.) et Xiao (J.). – Unified group fairness on federated learning, 2021.
22. Zhou (P.), Xu (H.), Lee (L. H.), Fang (P.) et Hui (P.). – Are You Left Out? An Efficient and Fair Federated Learning for Personalized Profiles on Wearable Devices of Inferior Networking Conditions. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 6, n2, 2022, pp. 1–25.