
The Effectiveness of Personalised Federated Learning Methods on Variations of Non-IID Data

Euodia Dodd

Department of Computer Science
University of Cambridge
ed581@cam.ac.uk

Abstract

1 Despite the recent advances in federated learning, the problem of dealing with het-
2 erogeneous data is still a substantial hurdle. In recent years, personalised federated
3 learning methods have been developed to capitalise on the diversity of client data
4 by creating more personalised models for individual clients. However, how these
5 methods cope with varying types of heterogeneity has yet to be investigated. In
6 this work, we evaluate the performance of state-of-the-art personalised federated
7 learning algorithms on varying types on non-IID data.

1 Introduction

In recent years, federated learning (FL) has gained interest due to its ability to leverage many different devices during the training process. The computational power of mobile devices has grown significantly allowing mobile devices with limited resources train machine learning models. Previously, edge devices would transmit data to a centralised server for training. However, this method suffers from privacy problems. With federated learning, each client trains their version of the model on local data and periodically communicates with a server which aggregates updates from the clients to train a centralised model. This helps preserve privacy by communication parameter updates instead of user data. In an ideal case, the global model obtained from aggregating client models should outperform local models due to more data being used. Whilst this has many benefits, it still suffers from some issues pertaining to privacy, fairness and how to deal with heterogeneous data. The aim of traditional FL algorithms such as FedAvg is to train a centralised model that performs well on average. We define f as the loss corresponding to user i .

$$\min_{w \in \mathbb{R}} f(w) := \frac{1}{n} \sum_{i=1}^n f_i(w)$$

For a supervised task with features x and labels y , we aim to solve:

$$f_i(w) := E_{(x,y) \sim p_i} [l_i(w; x, y)]$$

Due to the diversity of the clients available, it can not be assumed that the data for each client comes from the same distribution. Training a global model with non-independent and identically distributed (non-IID) data may lead to problems such as model parameter divergence and non-guaranteed convergence as different models may provide opposing gradient updates. It is widely known that an increase in statistical diversity leads to a significant increase in generalization error on local client data. Despite this challenge, non-IID data also provides an opportunity to create more personalised FL models for individual clients.

1.1 Problem formulation

It is generally well understood what non-IID means. But data can be non-IID in a variety of ways. Client data may be heterogeneous in different ways pertaining to the distribution of samples or labels across clients. These differences in the types of heterogeneity impacts the success of personalisation methods. We draw an example $(x, y) \sim P_i(x, y)$ from a client's local data. Heterogeneous data is defined as when $P_i \neq P_j$ for clients i and j . Kairouz et al. survey ways in which client distributions may be non-IID [1]:

Feature distribution skew - The marginal distributions $P_i(x)$ vary across clients even with shared $P(x | y)$.

Label distribution skew - Typically when clients are based in different geographical locations. For example, it may be more common to see snow in one location than another. $P_i(y)$ may differ across clients even with shared $P(x | y)$.

Concept drift (same label, different features or vice-versa) - When the same label may correspond to different features or the same features correspond to different labels. Could be due to cultural differences, weather effects etc. For example, road signs could look different in different parts of the world. $P_i(x | y)$ varies across clients even with shared $P(y)$ or $P_i(y | x)$ varies across clients even with shared $P(x)$.

Quantity skew - When certain clients have significantly more data than others

52 Despite the fact that real-world data typically contains a combination of these effects, most current
 53 works on personalised FL focus only on label distribution skew or quantity skew, i.e when the data
 54 comes from an existing, flat dataset which is then partitioned across clients by label and/or quantity.
 55 To our knowledge little work has been done to study these effects in isolation in an effort to understand
 56 their impact on the effectiveness of personalisation methods.

57 In this work, we investigate the effectiveness of state-of-the-art personalisation methods on different
 58 types of non-IID data. We focus on accuracy, convergence speed and stability of training. We
 59 hope that by understanding the impact different types of heterogeneity have on the performance of
 60 personalised FL algorithms, more methods can be developed to mitigate their effects. We evaluate
 61 the performance of these algorithms on non-IID variations of the MNIST dataset.

62 **2 Related work**

63 **2.1 Federated learning**

64 Federated stochastic gradient descent (FedSGD), typically regarded as a baseline, selects a fraction
 65 C of clients from which to sample gradient updates which are then averaged [2].

66 One of the first and most widely used FL algorithm is Federated Averaging (FedAvg) [3]. Federated
 67 averaging is an improvement on FedSGD where the local models are started from the same random
 68 initialisation and exchange model updates rather than gradient updates. There have been many
 69 proposed approaches to achieving personalisation in recent years. The main objective of these is to
 70 outperform FedAvg for local clients.

71 **2.2 Personalised Federated Learning**

72 Personalised FL has been studied widely with approaches including mixture-of-experts [4] and
 73 multi-task learning[5]. In this work we focus on two algorithms.

74 The first of these is pFedMe [6], an algorithm which uses Moreau envelopes as the clients' regularized
 75 loss function with the following l_2 -norm for each client:

$$f_i(\theta_i) + \frac{\lambda}{2} \|\theta_i - w\|^2$$

76 Where θ represents the personalised model for each client i and λ controls the weight of the global
 77 model w on the personalised model. In cases where client data is more abundant and useful, it
 78 may be more useful to appropriate a smaller significance to the global model. The personalised
 79 model can then optimise on local client data without straying too far from the global model.
 80 One advantage of pFedMe is the ability to update the global model similar to FedAvg, whilst
 81 optimising a personal model in parallel with low complexity. This can be split into a two step problem:

$$83 \min_{w \in \mathbb{R}^d} \{F(w) := \frac{1}{n} \sum_{i=1}^n F_i(w)\}, \text{ where } F_i(w) = \min_{\theta_i \in \mathbb{R}^d} \{f_i(\theta_i) + \frac{\lambda}{2} \|\theta_i - w\|^2\}$$

85 The optimised personalised model is then defined as the unique solution to the second part of the
 86 problem.

87 The results for pFedMe also demonstrate that it can achieve SOTA convergence speedup rate with
 88 carefully selected values of λ . Whilst larger values of λ increase convergence speed, it may also lead
 89 to divergence if too large.

90 Per-FedAvg, the second algorithm, was developed using the underlying idea of Model Agnostic Meta
 91 Learning (MAML) which assumes a limited budget to update a model when new data is observed
 92 [7]. It aims to find a global model to use as an initialisation which still performs well after each
 93 client performs a final gradient update step to optimize on it's own loss function. Per-FedAvg uses a
 94 similar regularising loss function to pFedMe but only optimises the first order approximation of the

loss function whereas pFedMe optimises this directly. However, Per-FedAvg requires computing or approximating the Hessian matrix which is computationally costly.

In the comparison of these algorithms conducted by Dinh et al., the heterogeneous data contained a combination of label skew and quantity skew. The results obtained show that pFedMe outperforms vanilla FedAvg and Per-FedAvg. Whether these results hold for other types of heterogeneity is an interesting problem.

3 Experiments

The main goal of this work is to evaluate the effectiveness of two SOTA personalised FL algorithms on variations of non-IID data. We focus on the MNIST dataset ¹ in particular. The three algorithms investigated are pFedMe, Per-FedAvg and FedAvg. In the case of pFedMe, we consider both the personalised and the global model.

Experimental Setup: Experimentation was conducted with a similar setup to the one used in pFedMe with the code adapted from the one used in the paper ². The MNIST dataset is partitioned amongst $N = 20$ clients. A learning rate of 0.005 was used with a subset of 5 clients. For the personalised version of pFedMe, a local learning rate of 0.1 was used. Larger values of K were found to show little improvement on the convergence of both the personalised and global model. As a result, a value of $K = 5$ was chosen. We set λ to 15. datasets are split with 75% training and 25% test data. Training was run for 20 local epochs with 800 rounds. An augmented version of the code used by Dinh et al. in the original paper was utilised. A l_2 -regularised multinomial logistic regression model, the same as the original paper, was also used. Each experiment was repeated 5 times and average to obtain the results.

Baseline: A baseline was established by training on IID data from the MNIST dataset. To generate non-IID data, the MNIST dataset was partitioned into equal parts and distributed amongst clients. This ensure no label skew, and the pre-shuffled data means the effects of label skew are negligible.

Heterogeneous data: In this work, label skew, quantity skew and concept drift are assessed separately. Quantity skew was achieved by partitioning the dataset along randomly generated split points. This created a local data size ranging between 216 and 10257 samples. As the dataset was already shuffled, the local data for each client contained a random number of samples from each of the 10 classes. The random seed was set to 1 to ensure reproducibility. We introduce concept drift by asymmetrically altering a label with probability P . The label l_i is changed to label l_{i+1} unless $l = 0$ in which case it is changed to 0. Doing so means images with the same features will have different labels. This experiment was conducted with $P = 0.1, 0.25$ and 0.5 to investigate the effects of varying degrees of concept drift. Label distribution skew was achieved by randomly assigning each client 2 classes out of the available 10. However, this introduced an element of quantity skew as each class contained a different number of samples though this effect is negligible.

4 Discussion

Baseline: As expected, training with homogeneous data was very smooth with fast convergence. An accuracy of over 90% was achieved for all algorithms.

¹<http://yann.lecun.com/exdb/mnist/>

²<https://github.com/CharlieDinh/pFedMe>

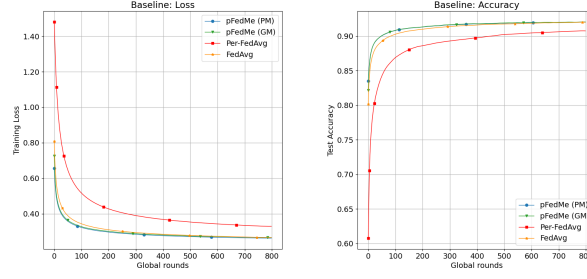


Figure 1: Accuracy and loss over 800 rounds for IID baseline data

133 **Quantity skew:** Quantity skew had very little impact. Over 91% accuracy was achieved for all
 134 algorithms with smooth training curves and fast convergence as shown in Figure 2. This is comparable
 135 to the results obtained with homogeneous data in the baseline experiment. This is likely due to the
 136 fact that the algorithms were developed and trained on datasets which contained large amounts of
 137 quantity skew. Intuitively, quantity skew is the most likely to occur and therefore the most widely
 138 tested form of heterogeneity. This may explain why the personalised FL algorithms are more robust
 139 to it.

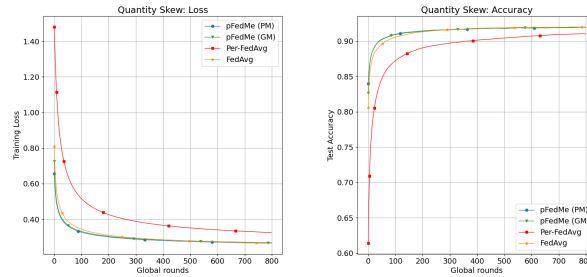


Figure 2: Accuracy and loss over 800 rounds for quantity skew

140 **Label skew:** The effects of label skew appear to be quite significant. The accuracy achieved by all
 141 algorithms was much lower than with IID data. However, Per-FedAvg and FedAvg converged faster
 142 and appeared much more stable than the two variations of pFedMe. To explain this, we take the
 143 example of $client_i$ with samples from classes 0,1. Assuming none of the other clients have samples
 144 from class 0 or 1, their models would be of much less relevance to $client_i$. This suggests that less
 145 weight should be applied to the global model and therefore a smaller value of λ should be selected to
 146 increase the accuracy of pFedMe. To test this hypothesis, we conducted further experiments with 2
 147 other values of lambda. The results show that a value of 15 outperformed both 10 and 20 suggesting
 148 that it is actually the optimal value.

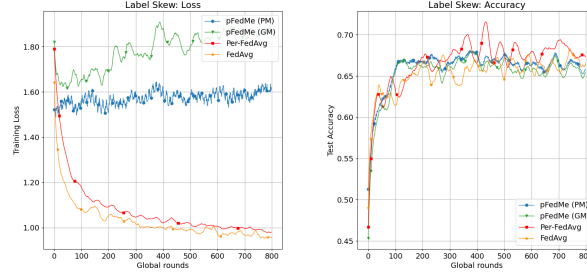


Figure 3: Accuracy and loss over 800 rounds for label skew

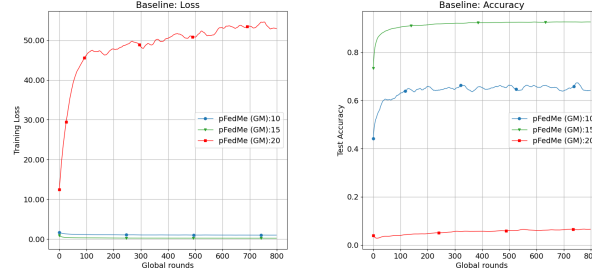


Figure 4: Accuracy and loss over 800 rounds for label skew with $\lambda = 10, 15, 20$

149 **Concept drift:** In the case of concept drift with $P = 0.1$, accuracy was also impacted slightly though
150 the results were not statistically significant. The effects of concept drift become more apparent at
151 $P = 0.25$ and 0.5 . At $P = 0.25$ the performance of all algorithms begin to suffer, dropping from
152 by almost 20% when compared to the baseline experiment and almost 16% compared to $P = 0.1$.
153 Vanilla FedAvg achieved a slightly higher accuracy than the personalised algorithms. As the concept
154 drift in this case was within local datasets, a higher value of lambda may have improved performance
155 as the input from other clients may have averaged out the effects of the corrupted labels.

156 An interesting effect is the decrease in accuracy over time. The maximum accuracy for each algorithm
157 is achieved within the first 100 or so rounds before the accuracy begins to decrease though this is only
158 observed for $P = 0.25$ and for Per-FedAvg with $P = 0.5$. It is unclear why this behaviour occurs but
159 a reduced learning rate or introducing dropout may go some way towards mitigating this. Accuracy
160 across all algorithms dropped again when $P = 0.5$ by over 40% compared to the baseline. Training
161 was also much less smooth.

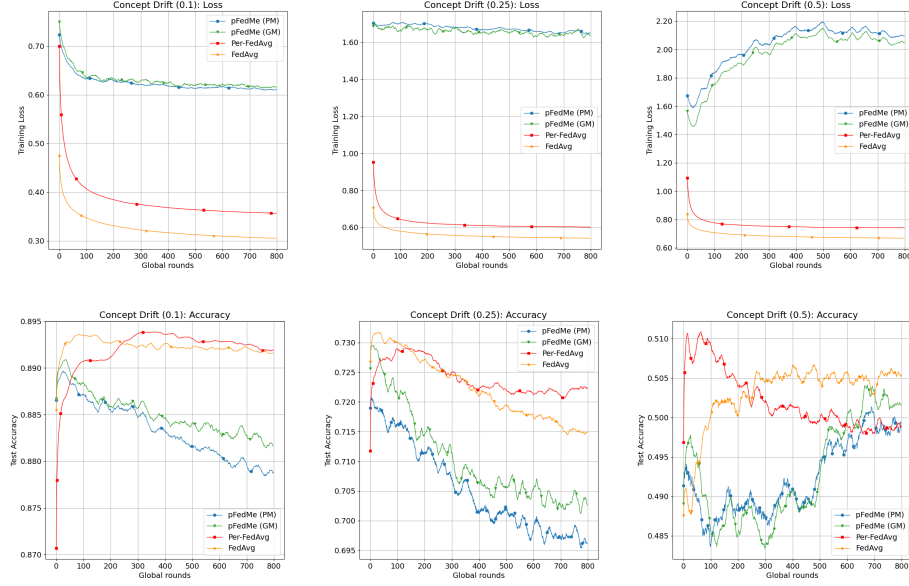


Figure 5: Accuracy and loss over 800 rounds for concept drift with $P = 0.1, 0.25$ and 0.5

Table 1: Max test accuracy obtained for each variation of non-IID data for each algorithm.

	Concept drift					
	Baseline(IID)	Quantity skew	Label skew	0.1	0.25	0.5
pFedMe(GM)	92.08	92.00	68.93	89.15	73.22	50.92
pFedMe(PM)	92.07	92.02	72.44	89.07	72.46	51.01
Per-FedAvg	90.77	91.08	76.00	89.44	73.30	51.57
FedAvg	92.02	92.03	70.33	89.40	73.34	51.16

5 Further work

Due to the long training time for pFedMe (approximately 90 minutes on an NVIDIA Tesla K80 GPU) some alterations had to be made to the methodology. In this work, we repeated the experiments 5 times instead of the 10 which may have affected the reliability of the results. To improve this, the experiments could be repeated more times. Additionally, the hyperparameters used in this work were obtained from the original experiments but may not have been optimal for the augmented datasets. Experimenting with different hyperparameters may have yielded a better performance, particularly the value of λ used for pFedMe. This work can be extended to include other SOTA personalised FL algorithms such as FedFomo.

Finally, the results of this work can be used to develop methods that mitigate the effects of heterogeneity in personalised FL training algorithms.

6 Conclusion

In this work, we evaluate the performance of pFedMe and per-FedAvg on 3 types of heterogeneous data. The results show that concept drift and label skew impacted the performance of both algorithms more severely than quantity skew. In the case of concept drift, vanilla FedAvg actually outperformed the personalised algorithms. Additionally, the notable difference in the way that they affect training can also be observed. This demonstrates that it is important to apply a more considered approach to

179 mitigate the effects of each kind of heterogeneity, rather than treating all heterogeneous data with a
180 one-fits-all solution.

References

- [1] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. A. Bonawitz, Z. Charles, G. Cormode, R. Cummings, R. G. L. D'Oliveira, S. E. Rouayheb, D. Evans, J. Gardner, Z. Garrett, A. Gascón, B. Ghazi, P. B. Gibbons, M. Gruteser, Z. Harchaoui, C. He, L. He, Z. Huo, B. Hutchinson, J. Hsu, M. Jaggi, T. Javidi, G. Joshi, M. Khodak, J. Konečný, A. Korolova, F. Koushanfar, S. Koyejo, T. Lepoint, Y. Liu, P. Mittal, M. Mohri, R. Nock, A. Özgür, R. Pagh, M. Raykova, H. Qi, D. Ramage, R. Raskar, D. Song, W. Song, S. U. Stich, Z. Sun, A. T. Suresh, F. Tramèr, P. Vepakomma, J. Wang, L. Xiong, Z. Xu, Q. Yang, F. X. Yu, H. Yu, and S. Zhao, "Advances and open problems in federated learning," *CoRR*, vol. abs/1912.04977, 2019.
- [2] R. Shokri and V. Shmatikov, "Privacy-preserving deep learning," in *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security, CCS '15*, (New York, NY, USA), p. 1310–1321, Association for Computing Machinery, 2015.
- [3] H. B. McMahan, E. Moore, D. Ramage, and B. A. y Arcas, "Federated learning of deep networks using model averaging," *CoRR*, vol. abs/1602.05629, 2016.
- [4] D. W. Peterson, P. Kanani, and V. J. Marathe, "Private federated learning with domain adaptation," *CoRR*, vol. abs/1912.06733, 2019.
- [5] V. Smith, C. Chiang, M. Sanjabi, and A. Talwalkar, "Federated multi-task learning," *CoRR*, vol. abs/1705.10467, 2017.
- [6] Y. Mansour, M. Mohri, J. Ro, and A. T. Suresh, "Three approaches for personalization with applications to federated learning," *CoRR*, vol. abs/2002.10619, 2020.
- [7] A. Fallah, A. Mokhtari, and A. E. Ozdaglar, "Personalized federated learning: A meta-learning approach," *CoRR*, vol. abs/2002.07948, 2020.